

Correctness of an STM Haskell Implementation

Manfred Schmidt-Schauss and David Sabel

Goethe-University, Frankfurt, Germany

Technical Report Frank-50

Research group for Artificial Intelligence and Software Technology

Institut für Informatik,

Fachbereich Informatik und Mathematik,

Johann Wolfgang Goethe-Universität,

Postfach 11 19 32, D-60054 Frankfurt, Germany

March 27, 2013

Abstract. A concurrent implementation of software transactional memory in Concurrent Haskell using a call-by-need functional language with processes and futures is given. The description of the small-step operational semantics is precise and explicit, and employs an early abort of conflicting transactions. A proof of correctness of the implementation is given for a contextual semantics with may- and should-convergence. This implies that our implementation is a correct evaluator for an abstract specification equipped with a big-step semantics.

1 Introduction

Due to the recent development in hardware and software, concurrent and parallel programming is getting more and more important, and thus there is a need for semantical investigations in concurrent programming. In this paper we investigate programming models of concurrent processes (threads) where several processes may access shared memory. Concurrent access and independency of processes lead to the well-known problems of conflicting memory use and thus requires protection of critical sections to ensure mutual exclusion. Over the years several programming primitives like locks, semaphores, monitors, etc. have been introduced and used to ensure this atomicity of memory operations in a concurrent setting. However, the explicit use of locking mechanisms is error-prone – the programmer may omit to set or release a lock resulting in deadlocks or race conditions – and it is also often inefficient, since setting too many locks may sequentialize program execution and prohibit concurrent evaluation. Another obstacle of lock-based concurrency is that composing larger programs from smaller ones is usually impossible [HMPH05].

A recent approach to overcome these difficulties is software transactional memory (STM) [ST95,ST97,HMPH05,HMPH08] where operations on the shared memory are viewed as *transactions* on the memory, i.e. several small operations (like read and write of memory locations) are compound to a transaction and then the system (the transaction manager) ensures that the transaction is performed in an atomic, isolated, and consistent manner. I.e., the programmer can assume that all the transactions are performed sequentially and isolated, while the runtime system may perform them concurrently and interleaved and may even try sequences of actions several times (invisible to the programmer). This removes the burden from the programmer to set locks and to keep track of them. Composability of transactions is another very helpful feature of STM. We focus on the Haskell-based approach in [HMPH05,HMPH08], which provides even more advantages like separation of different kinds of side effects (IO and STM) by monadic programming and the type system of Haskell. Memory-transactions are indicated in the program by marking a sequence of actions as atomic.

Though STM is easy to use for the programmer, implementing a *correct* transaction manager is considerably harder. Hurdles are an exact specification of the correctness properties, of course to provide an implementation, and to prove its correctness and the validity of the properties. In this paper we will address and solve these problems for a functional, concurrent language which models Concurrent Haskell extended by memory transactions. We start with a language and calculus *SHF* (STM-Haskell with futures) that is rather close to STM-Haskell and [HMPH05,HMPH08] and shares some ideas with the concurrent call-by-need process calculus investigated in [SSS11]. In contrast to STM-Haskell threads are

modelled by futures (which are however implementable in Concurrent Haskell using `unsafeInterleaveIO` in a safe manner [SSS12]), *SHF* has no explicit exceptions, it uses call-by-need evaluation and bindings in the environment for sharing, and – for simplicity – it is equipped with a monomorphic type system. We keep the monadic modelling for the separation of IO and STM and the composability, in particular the selection composability by `orElse`. A big-step operational semantics is given for *SHF*, which is obviously correct, since transactions are executed in isolation, and their effects on the shared memory are only observable after a successful execution. However, this semantics is *not implementable*, due to undecidable execution conditions in the big-step semantics. Thus the purpose of defining *SHF* standard reduction is not the implementation but the *specification* of a correct interpreter for *SHF*.

Secondly, we define a concurrent implementation of *SHF* by introducing the concurrent calculus *CSHF*. *CSHF* is close to a real implementation, since its operational semantics is formulated as a detailed, precise and complete small-step reduction semantics where all precautions and retries of the transactions are explicitly represented. I.e., the small-step reductions are defined with an appropriate granularity to make them implementable. Features of *CSHF* are registration of threads at transactional variables (TVars) and forced aborts of transactions in case of conflicts. All applicability conditions for the reductions are decidable. *CSHF* is designed to enable concurrent (i.e. interleaved) execution of threads as much as possible, i.e. there are only very short phases where internal locking mechanisms are used to prevent race conditions.

We also implemented *CSHF* as a prototypical library in Haskell¹, which provides the same core interface as Haskell’s STM implementation², and is (in contrast to *CSHF*) polymorphically typed. Internally it performs the steps as they are defined by the small-step semantics of *CSHF*. Our implementation behaves as expected and thus gives evidence of the correct design of *CSHF*.

The main goal of our investigation is to show that the concurrent implementation fulfills the specification. Here we use the strong criterion of contextual equivalence w.r.t. may- as well as should-convergence in both calculi, where may-convergence means that a process has the ability to evaluate to a successful process, and should-convergence means that the ability to become successful is never lost on every reduction path. Observing also should-convergence for contextual equivalence ensures that it is not permitted to transform a program P that cannot reach an error-state (i.e. a state that is not may-convergent) into a program that can reach an error-state. Compared with must-convergence, which requires termination on every computation path, should-convergence in combination with may-convergence judges busy-wait as non-erroneous behavior, and also contextual equivalence is invariant under restricting reductions to fair ones (see e.g. [RV07,SSS11]). Compared with definitions of contextual equivalences that include evaluation traces, using may- and should-convergence provides a coarser definition of equivalence.

In Main Theorem 4.3 we obtain the important result that the implementation mapping ψ is *observationally correct* [SSNSS08], i.e. for every context D and process P in *SHF* it holds: $D[P]$ may-converges (or should-converges, resp.), if, and only if $\psi(D)[\psi(P)]$ is may-convergent (should-convergent) in *CSHF*. Observational correctness thus shows that may- and should-convergence is preserved and reflected by the implementation ψ (i.e. ψ is convergence equivalent) and the tests used for contextual equivalence are translated in a compositional way. A direct consequence is also that ψ is *adequate*, i.e. the implementation preserves all contextual inequalities and thus does not introduce new equalities (which would introduce confusion in the implementation). Note that even the proof of convergence equivalence of ψ is a strong result, since it implies that *CSHF* is a correct evaluator for *SHF*.

Our notion of correctness is rather strong, since it implies the properties of the concurrent program that are aimed at in other papers like atomicity, opacity, asf. A surprising insight is that early abortion of conflicting transactions is *necessary* to make the implementation correct, where “early” means that a committing transaction must abort other conflicting transactions. This equivalence of *SHF* and *CSHF* has the following consequences: There is a guarantee that error-free *SHF*-programs are mapped to error-free *CSHF*-ones; error-freeness of a concurrent program means that no execution possibility leads to a state which is a dead end insofar as success can not be reached from this state. In particular if the *SHF*-program is free of deadlocks, then the translated *CSHF*-program is also free of deadlocks. The implementation will also guarantee a form of progress or liveness: If several transactions are active, and in the big-step variant there is at least one transaction that ends successfully, then also in the concurrent implementation at least one will terminate with success.

Comparing our results with the implementation of STM-Haskell in [HMPH08], our work is a justification for the correctness of most of the design decisions of their implementation. However, the imple-

¹ The documented source code and some instructions are available from <http://www.ki.cs.uni-frankfurt.de/research/stm>

² <http://hackage.haskell.org/package/stm>

$P, P_i \in Proc ::= P_1 \mid P_2 \mid \nu x.P \mid \langle ux \rangle \leftarrow e \mid x = e \mid \mathbf{xte}$ $e, e_i \in Expr ::= x \mid m \mid \lambda x.e \mid (e_1 \ e_2) \mid (c \ e_1 \ \dots \ e_{\text{ar}(c)})$ $\mid \mathbf{seq} \ e_1 \ e_2 \mid \mathbf{letrec} \ x_1 = e_1, \dots, x_n = e_n \ \mathbf{in} \ e$ $\mid \mathbf{case}_T \ e \ \mathbf{of} \ \mathit{alt}_{T,1} \ \dots \ \mathit{alt}_{T, T }$ <p style="text-align: center;">where $\mathit{alt}_{T,i} = (c_{T,i} \ x_1 \ \dots \ x_{\text{ar}(c_{T,i})} \rightarrow e_i)$</p> $m \in MExpr ::= \mathbf{return}_{\text{IO}} \ e \mid e_1 \gg_{\text{IO}} e_2 \mid \mathbf{future} \ e$ $\mid \mathbf{atomically} \ e \mid \mathbf{return}_{\text{STM}} \ e \mid e_1 \gg_{\text{STM}} e_2 \mid \mathbf{retry}$ $\mid \mathbf{orElse} \ e_1 \ e_2 \mid \mathbf{newTVar} \ e \mid \mathbf{readTVar} \ e \mid \mathbf{writeTVar} \ e$ $\tau, \tau_i \in Typ ::= \text{IO } \tau \mid \text{STM } \tau \mid \text{TVar } \tau \mid (T \ \tau_1 \ \dots \ \tau_n) \mid \tau_1 \rightarrow \tau_2$	<p><i>Functional values:</i> abstractions $\lambda x.e$ and constructor applications $(c \ e_1 \ \dots \ e_n)$</p> <p><i>Monadic values:</i> all monadic expressions $m \in MExpr$</p> <p><i>Values:</i> functional values and monadic values</p> <p><i>cx-values:</i> $(c \ x_1 \ \dots \ x_n)$ and monadic values where all arguments are variables</p> <p>(b) Functional values, monadic values, cx-values, and values</p> $(P_1 \mid P_2) \mid P_3 \equiv P_1 \mid (P_2 \mid P_3) \quad P_1 \mid P_2 \equiv P_2 \mid P_1$ $\nu x_1. \nu x_2. P \equiv \nu x_2. \nu x_1. P \quad P_1 \equiv P_2 \ \text{if} \ P_1 =_\alpha P_2$ $(\nu x. P_1) \mid P_2 \equiv \nu x. (P_1 \mid P_2) \ \text{if} \ x \notin FV(P_2)$ <p>(c) Structural Congruence</p>
(a) Syntax of Processes, Expressions, Monadic Expressions and Types	
$\mathbb{D} \in PC \quad ::= [\] \mid \mathbb{D} \mid P \mid P \mid \mathbb{D} \mid \nu x. \mathbb{D}$ $\mathbb{E} \in EC \quad ::= [\] \mid (\mathbb{E} \ e) \mid (\mathbf{case} \ \mathbb{E} \ \mathbf{of} \ \mathit{alts}) \mid (\mathbf{seq} \ \mathbb{E} \ e)$ $\widehat{M}_{\text{STM}} \in MC_{\text{STM}} ::= M_{\text{IO}}[\mathbf{atomically} \ \widehat{M}_{\text{STM}}]$ $\mathbb{L}_Q \in LC_Q \quad ::= \langle ux \rangle \leftarrow M_Q[\mathbb{F}] \mid \langle ux \rangle \leftarrow M_Q[\mathbb{F}[x_n]] \mid x_n = \mathbb{E}_n[x_{n-1}] \mid \dots \mid x_2 = \mathbb{E}_2[x_1] \mid x_1 = \mathbb{E}_1$ <p style="text-align: center;">where $Q \in \{\text{IO}, \text{STM}\}$ and $\mathbb{E}_1 \neq [\]$</p>	$\mathbb{F} \in FC \quad ::= \mathbb{E} \mid (\mathbf{readTVar} \ \mathbb{E}) \mid (\mathbf{writeTVar} \ \mathbb{E} \ e)$ $M_{\text{IO}} \in MC_{\text{IO}} ::= [\] \mid M_{\text{IO}} \gg_{\text{IO}} e$ $\widehat{M}_{\text{STM}} \in \widehat{MC}_{\text{STM}} ::= [\] \mid \widehat{M}_{\text{STM}} \gg_{\text{STM}} e \mid \mathbf{orElse} \ \widehat{M}_{\text{STM}} \ e$
(d) Process-, Monadic-, Evaluation-, and Forcing-Contexts	
$\mathbf{future} \quad ::= (\text{IO } \alpha) \rightarrow \text{IO } \alpha \quad \mathbf{readTVar} \quad ::= (\text{TVar } \alpha) \rightarrow \text{STM } \alpha$ $\mathbf{atomically} \quad ::= (\text{STM } \alpha) \rightarrow \text{IO } \alpha \quad \mathbf{writeTVar} \quad ::= (\text{TVar } \alpha) \rightarrow \alpha \rightarrow \text{STM } ()$ $\mathbf{return}_{\text{STM}} \quad ::= \alpha \rightarrow \text{STM } \alpha \quad \mathbf{newTVar} \quad ::= \alpha \rightarrow \text{STM}(\text{TVar } \alpha)$ $\mathbf{return}_{\text{IO}} \quad ::= \alpha \rightarrow \text{IO } \alpha \quad \mathbf{retry} \quad ::= (\text{STM } \alpha)$	$\gg_{\text{IO}} \quad ::= (\text{IO } \alpha_1) \rightarrow (\alpha_1 \rightarrow \text{IO } \alpha_2) \rightarrow \text{IO } \alpha_2$ $\gg_{\text{STM}} \quad ::= (\text{STM } \alpha_1) \rightarrow (\alpha_1 \rightarrow \text{STM } \alpha_2) \rightarrow \text{STM } \alpha_2$ $\mathbf{orElse} \quad ::= (\text{STM } \alpha) \rightarrow (\text{STM } \alpha) \rightarrow (\text{STM } \alpha)$
(e) Polymorphic types of operators	
$\frac{\Gamma(x) = \tau, \quad \Gamma \vdash e :: \text{IO } \tau}{\Gamma \vdash \langle ux \rangle \leftarrow e :: \mathbf{wt}} \quad \frac{\Gamma(x) = \tau, \quad \Gamma \vdash e :: \tau}{\Gamma \vdash x = e :: \mathbf{wt}} \quad \frac{\Gamma \vdash P_1 :: \mathbf{wt}, \quad \Gamma \vdash P_2 :: \mathbf{wt}}{\Gamma \vdash P_1 \mid P_2 :: \mathbf{wt}} \quad \frac{\Gamma(x) = \text{TVar } \tau, \quad \Gamma \vdash e :: \tau}{\Gamma \vdash \mathbf{xte} :: \mathbf{wt}} \quad \frac{\Gamma \vdash P :: \mathbf{wt}}{\Gamma \vdash \nu x.P :: \mathbf{wt}}$ $\frac{\Gamma(x) = \tau}{\Gamma \vdash x :: \tau} \quad \frac{\Gamma(x) = \tau_1, \quad \Gamma \vdash e :: \tau_2}{\Gamma \vdash (\lambda x.e) :: \tau_1 \rightarrow \tau_2} \quad \frac{\Gamma \vdash e_1 :: \tau_1 \rightarrow \tau_2, \quad \Gamma \vdash e_2 :: \tau_1}{\Gamma \vdash (e_1 \ e_2) :: \tau_2} \quad \frac{\Gamma \vdash e_i :: \tau_i, i = 1, 2 \quad \tau_1 = \tau_3 \rightarrow \tau_4 \ \text{or} \ \tau_1 = (T \ \dots)}{\Gamma \vdash (\mathbf{seq} \ e_1 \ e_2) :: \tau_2}$ $\frac{\forall i : \Gamma(x_i) = \tau_i, \quad \forall i : \Gamma \vdash e_i :: \tau_i, \quad \Gamma \vdash e :: \tau}{\Gamma \vdash (\mathbf{letrec} \ x_1 = e_1, \dots, x_n = e_n \ \mathbf{in} \ e) :: \tau} \quad \frac{\Gamma \vdash e :: \tau_1 \ \text{and} \ \tau_1 = (T \ \dots), \quad \forall i : \Gamma \vdash (c_i \ x_{1,i} \ \dots \ x_{n_i,i}) :: \tau_i, \quad \forall i : \Gamma \vdash e_i :: \tau_2}{\Gamma \vdash (\mathbf{case}_T \ e \ \mathbf{of} \ (c_1 \ x_{1,1} \ \dots \ x_{n_1,1} \rightarrow e_1) \ \dots \ (c_m \ x_{1,m} \ \dots \ x_{n_m,m} \rightarrow e_m)) :: \tau_2}$ $\frac{\forall i : \Gamma \vdash e_i :: \tau_i, \quad \tau_1 \rightarrow \dots \rightarrow \tau_n \rightarrow \tau_{n+1} \in \mathbf{types}(c)}{\Gamma \vdash (c \ e_1 \ \dots \ e_{\text{ar}(c)}) :: \tau_{n+1}} \quad \text{where } c \text{ is a constructor, or a monadic operator and } \mathbf{types}(c) \text{ is the set of monomorphic types of a constructor or a monadic operator } c$	
(f) Typing rules	

Fig. 1. The calculus SHF, syntax and contexts

mentations behave slightly different which will be discussed in Section 5. The reason for not using the STM-Haskell implementation in [HMPH08] is that “formalizing” this implementation as a calculus is not easily possible, since e.g. it requires a test for pointer equality, which is cumbersome in a call-by-need calculus with indirections (represented by binding chains).

Outline. In Section 2 we define the (specification) calculus *SHF*, by explaining its syntax and operational semantics. The concurrent implementation in form of the calculus *CSHF* is introduced in Section 3. In Section 4 we provide the translation $\psi : SHF \rightarrow CSHF$ and prove that ψ is observationally correct which shows correctness of the implementation. We discuss related work in Section 5 and finally conclude in Section 6.

2 The SHF-Calculus

The syntax of *SHF* and its *processes Proc* is in Fig. 1(a). We assume a countable infinite set of variables (denoted by x, y, z, \dots) and a countable infinite set of *identifiers* (denoted by u, u') to identify threads.

In a *parallel composition* $P_1 \mid P_2$ the processes P_1 and P_2 run concurrently, and the *name restriction* $\nu x.P$ restricts the scope of x to process P . A *concurrent thread* $\langle ux \rangle \leftarrow e$ has *identifier* u and evaluates the expression e binding the result to the variable x (called the *future* x). A process has usually one

distinguished thread, the *main thread*, denoted by $\langle u|x \rangle \stackrel{\text{main}}{\longleftarrow} e$. Evaluation of the main thread is enforced, which does not hold for other threads. *TVars* $x \mathbf{t} e$ are the transactional (mutable) variables with name x and content e . They can only be accessed inside STM-transactions. *Bindings* $x = e$ model globally shared expressions, where we call x a *binding variable*. Futures and names of TVars are called *component-names*.

We assume a partitioned set of *data constructors*, where each family represents a type constructor T . The T -data constructors are ordered (denoted with $c_1, \dots, c_{|T|}$). Each type constructor T and each data constructor c has an arity $\text{ar}(T) \geq 0$ ($\text{ar}(c) \geq 0$, resp.). The functional language has variables, *abstractions* $\lambda x.e$, and *applications* $(e_1 e_2)$, *constructor applications* $(c e_1 \dots e_{\text{ar}(c)})$, *case-expressions* where for every type constructor T there is one case_T -construct. For case_T there is exactly one *case-alternative* $(c_{T,i} x_1 \dots x_{\text{ar}(c_{T,i})} \rightarrow e_i)$ for every constructor $c_{T,i}$ of type constructor T where the variables x_i become bound with scope e_i . In a *pattern* $c_{T,i} x_1 \dots x_{\text{ar}(c_{T,i})}$ the variables x_i must be pairwise distinct. We sometimes abbreviate the *case-alternatives* with *alts*. Further constructs are *seq-expressions* $(\text{seq } e_1 e_2)$ for strict evaluation, and *letrec-expressions* to implement local sharing and recursive bindings. In *letrec* $x_1 = e_1, \dots, x_n = e_n$ **in** e the variables x_i must be pairwise distinct and the bindings $x_i = e_i$ are recursive, i.e. the scope of x_i is e_1, \dots, e_n and e .

Monadic expressions comprise variants for the IO- and the STM-monad of the “bind” operator $\gg=$ for sequential composition of actions, and the *return*-operator. For the STM-monad *newTVar*, *readTVar*, and *writeTVar* are available to create and access TVars, the primitive *retry* to abort and restart the STM-transaction, and *orElse* $e_1 e_2$ to compose an STM-transaction from e_1, e_2 : *orElse* returns if e_1 is successful and if it catches a *retry* in e_1 then it proceeds with e_2 . For the IO-monad the *future*-operator creates threads, and *atomically* lifts an STM-transaction into the IO-monad by executing the transaction.

Variable binders are introduced by abstractions, *letrec*, *case-alternatives*, and $\nu x.P$. This induces a notion of free and bound variables and α -renaming and α -equivalence (denoted by $=_\alpha$). Let $FV(P)$ ($FV(e)$, resp) be the free variables of process P (expression e , resp.). We assume the *distinct variable convention* to hold, and also assume α -renaming is implicitly performed, if necessary.

A *context* is a process or an expression with a hole (denoted by $[\cdot]$). We write $C[e]$ ($C[P]$, resp.) for filling the hole of context C by expression e (process P , resp.). For processes we use a *structural congruence* \equiv to equate obviously equal processes, which is the least congruence satisfying the equations of Fig. 1(c).

SHF is equipped with a monomorphic type system. The syntax of monomorphic types Typ is given in Fig. 1(a), where $(\text{IO } \tau)$ means a monadic IO-action with return type τ , $(\text{STM } \tau)$ means an STM-transaction action, $(\text{TVar } \tau)$ stands for a TVar-reference with content type τ , $(T \tau_1 \dots \tau_n)$ is a type for an n -ary type constructor T , and $\tau_1 \rightarrow \tau_2$ is a function type. Although the type system is monomorphic, we “overload” the constructors and the monadic operators by assuming that they have a polymorphic type according to the usual conventions, however, in the language they are used as monomorphic. The polymorphic types of the monadic operators are shown in Fig. 1(e) where α, α_i are type variables.

To fix the types during reduction and for transformations, we assume that every variable x is explicitly typed and thus has a built-in type $\Gamma(x)$. For contexts, we assume that the hole $[\cdot]$ is typed and carries a type label. The notation $\Gamma \vdash e :: \tau$ ($\Gamma \vdash P :: \mathbf{wt}$, resp.) means that expression e (process P , resp.) can be typed with type τ (can be typed, resp.) using Γ . The typing rules are given in Fig. 1(f). Note that $\langle u|x \rangle \leftarrow e$ is well-typed if x is of type τ and e is of type $\text{IO } \tau$, and that an STM- or an IO-type for the first argument of *seq* is forbidden, to enable that the monad laws hold (see [SSS12]).

Definition 2.1. *A variable x is an introduced variable if it is a binding variable or a component-name. A process is well-formed, if all introduced variables and identifiers are pairwise distinct, and it has at most one main thread.*

A process P is well-typed iff P is well-formed and $\Gamma \vdash P :: \mathbf{wt}$ holds. An expression e is well-typed with type τ iff $\Gamma \vdash e :: \tau$ holds.

We define a reduction relation $\xrightarrow{\text{SHF}}$ (called *standard reduction*) for *SHF*. This reduction relation is on the one hand a small-step reduction, since several $\xrightarrow{\text{SHF}}$ -steps are performed sequentially before evaluation stops. However, it uses as a condition the success of the intermediate reduction relation $\xrightarrow{\text{SHFA}}$ in a big-step manner for the atomic execution of an STM-transaction (i.e. the evaluation of e in *atomically* e and of *orElse* $e e'$).

Definition 2.2. *A well-formed process P is successful, if P has a main thread of the form $\langle u|x \rangle \stackrel{\text{main}}{\longleftarrow} \text{return } e$, i.e. $P \equiv \nu x_1. \dots \nu x_n. (\langle u|x \rangle \stackrel{\text{main}}{\longleftarrow} \text{return } e \mid P')$.*

Monadic STM Computations:

$$\begin{aligned}
(\text{lunit}_{\text{STM}}) \quad & \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{return}_{\text{STM}} e_1 \gg_{\text{STM}} e_2] \xrightarrow{\text{SHFA}} \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[e_2 e_1] \\
(\text{read}) \quad & \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{readTVar } x \mid \text{x t e}] \xrightarrow{\text{SHFA}} \nu z. (\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{return}_{\text{STM}} z] \mid z = e \mid \text{x t z}) \\
(\text{write}) \quad & \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{writeTVar } x e_2 \mid \text{x t e}_1] \xrightarrow{\text{SHFA}} \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{return}_{\text{STM}} ()] \mid \text{x t e}_2 \\
(\text{nvar}) \quad & \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{newTVar } e] \xrightarrow{\text{SHFA}} \nu x. (\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{return}_{\text{STM}} x] \mid \text{x t e}) \\
(\text{ortry}) \quad & \frac{\nu X. \mathbb{D}[\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{orElse } e_1 e_2]] \xrightarrow{\text{SHFA},*} \nu X. \mathbb{D}'[\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{orElse retry } e_2]]}{\nu X. \mathbb{D}[\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{orElse } e_1 e_2]] \xrightarrow{\text{SHFA}} \nu X. \mathbb{D}[\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[e_2]]} \\
(\text{orret}) \quad & \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{orElse } (\text{return}_{\text{STM}} e_1) e_2] \xrightarrow{\text{SHFA}} \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[e_1] \\
(\text{retryup}) \quad & \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{retry} \gg_{\text{STM}} e_1] \xrightarrow{\text{SHFA}} \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{STM}}[\text{retry}]
\end{aligned}$$

Monadic IO Computations:

$$\begin{aligned}
(\text{lunit}_{\text{IO}}) \quad & \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{return}_{\text{IO}} e_1 \gg_{\text{IO}} e_2] \xrightarrow{\text{SHF}} \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{IO}}[e_2 e_1] \\
(\text{fork}) \quad & \langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{future } e] \xrightarrow{\text{SHF}} \nu z, u'. (\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{return } z] \mid \langle u'z \rangle \Leftarrow e) \\
& \text{where } z, u' \text{ are fresh and the created thread is not the main thread} \\
(\text{unIO}) \quad & \langle \text{u}ly \rangle \Leftarrow \text{return } e \xrightarrow{\text{SHF}} y = e \quad \text{if the thread is not the main-thread} \\
(\text{atomic}) \quad & \frac{\nu X. \mathbb{D}_1[\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{atomically } e]] \xrightarrow{\text{SHFA},*} \nu X. \mathbb{D}'_1[\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{atomically } (\text{return}_{\text{STM}} e')]]}{\mathbb{D}[\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{atomically } e]] \xrightarrow{\text{SHF}} \mathbb{D}'[\langle \text{u}ly \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{return}_{\text{IO}} e']] } \\
& \text{where } \mathbb{D} = \nu X. (\mathbb{D}_1 \mid \mathbb{D}_2) \text{ and } \mathbb{D}_1 \text{ contains all bindings and TVars of } \mathbb{D}, \\
& \mathbb{D}_2 \text{ contains all futures of } \mathbb{D}, \text{ and } \mathbb{D}' = \nu X. (\mathbb{D}'_1 \mid \mathbb{D}_2)
\end{aligned}$$

Functional Evaluation:

$$\begin{aligned}
(\text{feval}_{\text{IO}}) \quad & P \xrightarrow{\text{SHF}} P', \text{ if } P \xrightarrow{a} P' \text{ for } a \in \{\text{cp}_{\text{IO}}, \text{abs}_{\text{IO}}, \text{mkb}_{\text{IO}}, \text{lbeta}_{\text{IO}}, \text{case}_{\text{IO}}, \text{seq}_{\text{IO}}\} \\
(\text{feval}_{\text{STM}}) \quad & P \xrightarrow{\text{SHFA}} P', \text{ if } P \xrightarrow{a} P' \text{ for } a \in \{\text{cp}_{\text{STM}}, \text{abs}_{\text{STM}}, \text{mkb}_{\text{STM}}, \text{lbeta}_{\text{STM}}, \text{case}_{\text{STM}}, \text{seq}_{\text{STM}}\}
\end{aligned}$$

The reductions with parameter $Q \in \{\text{STM}, \text{IO}\}$ are defined as follows:

$$\begin{aligned}
(\text{cp}_Q) \quad & \mathbb{L}_Q[x_1 \mid x_1 = x_2 \mid \dots \mid x_{n-1} = x_n \mid x_n = v] \rightarrow \mathbb{L}_Q[v \mid x_1 = x_2 \mid \dots \mid x_{n-1} = x_n \mid x_n = v], \\
& \text{if } v \text{ is an abstraction, a cx-value or a component-name} \\
(\text{abs}_Q) \quad & \mathbb{L}_Q[x_1 \mid x_1 = x_2 \mid \dots \mid x_{m-1} = x_m \mid x_m = c e_1 \dots e_n \\
& \rightarrow \nu y_1, \dots, y_n. (\mathbb{L}_Q[c y_1 \dots y_n] \mid x_1 = x_2 \mid \dots \mid x_{m-1} = x_m \mid x_m = c y_1 \dots y_n \mid y_1 = e_1 \mid \dots \mid y_n = e_n) \\
& \text{if } c \text{ is a constructor, or } \text{return}_{\text{STM}}, \text{return}_{\text{IO}}, \gg_{\text{STM}}, \gg_{\text{IO}}, \text{orElse}, \text{atomically}, \text{readTVar}, \text{writeTVar}, \text{newTVar}, \text{or} \\
& \text{future and } n \geq 1, \text{ and some } e_i \text{ is not a variable.} \\
(\text{mkb}_Q) \quad & \mathbb{L}_Q[\text{letrec } x_1 = e_1, \dots, x_n = e_n \text{ in } e] \rightarrow \nu x_1, \dots, x_n. (\mathbb{L}_Q[e] \mid x_1 = e_1 \mid \dots \mid x_n = e_n) \\
(\text{lbeta}_Q) \quad & \mathbb{L}_Q[(\lambda x. e_1) e_2] \rightarrow \nu x. (\mathbb{L}_Q[e_1] \mid x = e_2) \\
(\text{case}_Q) \quad & \mathbb{L}_Q[\text{case}_T (c e_1 \dots e_n) \text{ of } \dots ((c y_1 \dots y_n) \rightarrow e) \dots] \rightarrow \nu y_1, \dots, y_n. (\mathbb{L}_Q[e] \mid y_1 = e_1 \mid \dots \mid y_n = e_n), \text{ if } n > 0 \\
(\text{case}_Q) \quad & \mathbb{L}_Q[\text{case}_T c \text{ of } \dots (c \rightarrow e) \dots] \rightarrow \mathbb{L}_Q[e] \\
(\text{seq}_Q) \quad & \mathbb{L}_Q[(\text{seq } v e)] \rightarrow \mathbb{L}_Q[e], \quad \text{if } v \text{ is a functional value}
\end{aligned}$$

$$\text{Closure: } \frac{P_i \equiv D[P'_i], P'_1 \xrightarrow{\text{SHF}} P'_2}{P_1 \xrightarrow{\text{SHF}} P_2} \quad \frac{P_i \equiv D[P'_i], P'_1 \xrightarrow{\text{SHFA}} P'_2}{P_1 \xrightarrow{\text{SHFA}} P_2}$$

Fig. 2. The calculus *SHF*, reductions

Definition 2.3. *The standard reduction rules are given in Fig. 2 where the used contexts are defined in Fig. 1(d). Standard reductions are permitted only for well-formed and non-successful processes.*

In the following we will denote the transitive closure of a relation \xrightarrow{R} by $\xrightarrow{R,+}$ and the reflexive-transitive closure by $\xrightarrow{R,*}$.

We define the \xrightarrow{SHF} -redex: For (lunit), (fork), it is the expression in the context \mathbb{M} , for (unIO), it is $\langle u!y \rangle \leftarrow \mathbf{return} \ e$, for (mkb), (lbeta), (case), (seq), (cp), (abs) it is the expression (or variable) in the context \mathbb{L} . We define the \xrightarrow{SHFA} -redex: For (lunit_{STM}), (read), (write), (nvar), (ortry), (orret) it is the expression in the context \mathbb{M}_{STM} , for (mkb_{STM}), (lbeta_{STM}), (case_{STM}), (seq_{STM}), (cp_{STM}), (abs_{STM}) it is the expression (or variable) in the context \mathbb{L}_{STM} .

We explain the standard reduction rules of Fig. 2. The standard (big-step) reduction is \xrightarrow{SHF} , whereas $\xrightarrow{SHFA,*}$ defines the evaluation of STM-transactions. The $\xrightarrow{SHFA,*}$ -reduction sequences will only be performed, if they terminate successfully, see rule (atomic).

The rules are divided into three classes: Rules for monadic computations in the STM-monad, rules for monadic computations in the IO-monad, and rules for functional evaluation.

We start with explaining the most interesting rules – the rules for the STM-monad: The rule (lunit) implements the semantics of the monadic sequencing operator $\gg=$. The rules (read), (write), and (nvar) access and create TVars. The rule (ortry) is a big-step rule: if a $\xrightarrow{SHFA,*}$ -reduction sequence starting with $\mathbf{orElse} \ e_1 \ e_2$ ends in $\mathbf{orElse} \ \mathbf{retry} \ e_2$, then the effects are ignored, and $\mathbf{orElse} \ e_1 \ e_2$ is replaced by e_2 . If the reduction of e_1 ends with $\mathbf{return} \ e$, then rule (orret) is used to keep the result as the result of $\mathbf{orElse} \ e_1 \ e_2$.

For the IO-monad there is also a rule (lunit) for the bind-operator. The rule (atomic) executes an STM-action in the IO-monad. It is also a big-step rule. If for a single thread the $\xrightarrow{SHFA,*}$ -reduction successfully produces a return, then the transaction is performed in one step of the \xrightarrow{SHF} -reduction. If the $\xrightarrow{SHFA,*}$ -reduction ends in a \mathbf{retry} , in a stuck expressions or does not terminate, then there is no \xrightarrow{SHF} -reduction, and hence it is omitted from the operational semantics.

The rule (fork) spawns a new concurrent thread and returns the newly created future as the result. The rule (unIO) binds the result of a monadic computation to a functional binding, i.e. the value of a concurrent future becomes accessible.

The rules for functional evaluation implement call-by-need reduction. The rules (cp) and (abs) inline a needed binding $x = e$ where e must be an abstraction, a cx-value, or a component name. To implement call-by-need evaluation the arguments of constructor applications and monadic expressions are shared by new bindings, similar to lazy copying [vEPS90]. Since the variable (binding-) chains are transparent, there is no need to copy binding-variables to other places in the expressions. The rule (mkb) moves \mathbf{letrec} -bindings into the global bindings. The rule (lbeta) is the sharing variant of β -reduction. The (case)-reduction reduces a \mathbf{case} -expression, where perhaps bindings are created to implement sharing. The (seq)-rule evaluates a \mathbf{seq} -expression.

Since the reduction rules only introduce variables which are fresh and never introduce a main thread, \xrightarrow{SHF} preserves well-formedness. Also type preservation holds, since every redex keeps the type of subexpressions.

Contextual equivalence equates processes P_1, P_2 if their observable behavior is indistinguishable if P_1 and P_2 are plugged into any process context. For nondeterministic (and concurrent) calculi observing may-convergence, i.e. whether a process can be reduced to a successful process, is *not* sufficient and thus we will observe may-convergence and should-convergence (see [RV07,SSS08]).

Definition 2.4. *A process P may-converges (written as $P\Downarrow$), iff it is well-formed and reduces to a successful process, (see Definition 2.2), i.e. $P\Downarrow$ iff P is well-formed and $\exists P' : P \xrightarrow{SHF,*} P' \wedge P'$ successful. If $P\Downarrow$ does not hold, then P must-diverges written as $P\Uparrow$. A process P should-converges (written as $P\Downarrow$), iff it is well-formed and remains may-convergent under reduction, i.e. $P\Downarrow$ iff P is well-formed and $\forall P' : P \xrightarrow{SHF,*} P' \implies P'\Downarrow$. If P is not should-convergent then we say P may-diverges written as $P\Uparrow$, which is also equivalent to $\exists P' : P \xrightarrow{SHF,*} P' \wedge P'\Uparrow$.*

Contextual approximation \leq_c and contextual equivalence \sim_c on processes are defined as $\leq_c := \leq_\downarrow \cap \leq_\uparrow$ and $\sim_c := \leq_c \cap \geq_c$ where for $\zeta \in \{\downarrow, \uparrow\}$: $P_1 \leq_\zeta P_2$ iff $\forall \mathbb{D} \in PC : \mathbb{D}[P_1]\zeta \implies \mathbb{D}[P_2]\zeta$.

The definition of \xrightarrow{SHF} implies that non-wellformed processes are always must-divergent. Also, the process construction by $\mathbb{D}[P]$ is always well-typed if P is well-typed, since we assume that variables have a built-in type.

Clearly, SHF implements memory transactions in a correct way: Every single transaction is executed atomically and isolated from any other transaction. So SHF is a correct specification of STM. However, the semantics has two drawbacks:

- It is impossible to implement the standard reduction: The rules (atomic) and (ortry) have undecidable preconditions, since they include checking whether a \xrightarrow{SHFA} -reduction halts successfully.
- Since the transactions are executed sequentially, there is only poor concurrency.

In the next section we will give an implementation using SHF as specification, which overcomes those drawbacks.

3 A Concurrent Implementation of STM Evaluation

We introduce a concurrent (small-step) evaluation enabling much more concurrency which is closer to an abstract machine than the big-step semantics. The executability of every single step is decidable, and every state has a finite set of potential successors. The undecidable conditions in the rules (atomic) and (ortry) of SHF are now implemented by tentatively executing the transaction, under concurrency, thus allowing other threads to execute. Transaction execution should guarantee an equivalent linearized execution of transactions. Our ultimate goal is to show that the concurrent execution is *correct*, i.e. to show that it is semantically equivalent to the big-step reduction defined for SHF . However, our approach uses minimal locking times and thus has a potential danger of conflicts and retries. Instead of using an optimistic read/write approach which performs a rollback in case of a conflict [HLR10], we will follow a similar, but slightly more pessimistic read and write: there are no locks at the start of a transaction, initial reads copy contents to the local store, other reads and writes are local; only at the end of a successful transaction there is a commit phase with global writes such that updates become visible to other transactions where a real, but internal, locking is used for a short time. To have a correct overall execution, conflicting transactions will be stopped by sending them a retry-notification (so-called early conflict detection), where the knowledge of the potential conflicts is memorized at the TVars.

We view a transaction as a function $V_1 \rightarrow V_2$ from a set of (read) input TVars V_1 to a set of modified TVars V_2 where V_1, V_2 may have common elements. The guarantee must be that at the end of transaction execution, the complete transaction could be atomically and instantaneously executed on the TVars V_1, V_2 . During this commit phase other transactions that have read any variable of V_2 , but are not finished yet, need to be aborted (restarted) due to a conflict.

More concretely, the ideas of the concurrent implementation are as follows:

- For every transaction there is a *local copy* for every TVar it accesses. All reads and writes are performed in the local copies invisible to other concurrently executed threads.
- For later conflict detection every global TVar is equipped with a set of registered threads. Before obtaining a local copy the thread identifier is added to this set.
- For book-keeping of the accessed TVars every thread has a transaction log, where read, written, and created TVars are remembered.
- If an executing transaction (thread) results in **retry**, its local copies of the TVars are removed, the registrations on the TVars are removed and the transaction starts over again. If the **retry** occurs inside an **orElse** then the treatment is different (see next item).
- To implement the nesting of **orElse**, the local copies of the TVars are stacks, where the stack depth matches the depth of the **orElse**-nesting. In case that a transaction stops with **retry** inside the left argument of **orElse**, all these stacks are popped and the right argument of **orElse** is executed. However, read TVars must remain in the transaction log, since their values may have an influence on the control. For the same reason, also the registration at the global TVars must not be changed.
- If a transaction executed all its actions, then it starts its *commit phase*. To ensure consistency, first all read and all to-be-written TVars are locked by the thread identifier of the committing transaction. This locking is done in one atomic step. After the locks are set, no other thread is able to access the locked TVars. Then there are several steps performed which may be interleaved by other transactions (accessing other TVars):
 - The committing thread removes its registration from all TVars it has accessed.

- The committing thread reads the registration sets to obtain the thread identifiers of all the conflicting transactions which have to be aborted. The corresponding threads are *notified* by the committing transaction to abort and restart their transaction. This can be seen as throwing an exception from the committing thread to all other threads which have become inconsistent now. This is indeed implemented in this way in our prototype. In the operational semantics below, simply the current code of the other transactions is replaced by `retry` and thus notified threads perform as if they got a (transaction) `retry`: first they unregister their thread identifier and then restart.
- The committing thread updates the global content of the TVars and then removes the locks.
- The final steps of the committing thread are to add newly created TVars as global TVars and then to remove all the local copies.

Another design decision in the calculus description below is that sharing by bindings is carefully maximized, including transparent variable-variable-binding chains. One reason is that this leads to manageable correctness proofs. Another design principle is to avoid negative conditions: There must not be any conditions that rely on a test whether there are no occurrences of TVars or threads satisfying certain conditions.

Now we detail on this idea and introduce the calculus *CSHF* which has a concurrent evaluation of transactions and a modified *SHF*-syntax: there is an addition of labels, of memory-cells at threads and a stack for the `orElse`s. First we describe the syntax changes of the language, and then exactly describe the rules, which are close to a concurrent abstract machine for SHF.

Definition 3.1 (Syntax of CSHF). *The syntax of CSHF-processes is almost the same as for SHF-processes where, however, there are some extensions and changes. Instead of TVars $x\mathbf{t}e$ there are two constructs:*

1. Global TVars are represented by $x\mathbf{t}g\ e\ u\ g$, where the additional third argument u is a locking label and may be empty (written as $-$) or a thread identifier u that locks the TVar, and g is a list of thread identifiers for those threads that want to be notified for a retry, when x is updated.
2. A stack of thread-local TVars is represented by $u\mathbf{t}l\ s$, where u is a thread identifier and s is a stack of sets with local TVars $x\mathbf{t}l\ e$ as elements, where x is a name of a TVar and e is an expression.

A thread may have a transaction log, which is only available if a thread is within a transaction. It is written over the thread-arrow as $\langle u\mathbf{t}y \rangle \xleftarrow{T,L;K} e$ where T is a set of TVars (i.e. the names of the TVars) that are read during the transaction; L is a stack of triples (L_a, L_n, L_w) where L_a is a set of the names of all TVars which are accessed during the transaction, L_n a set of names of newly generated TVars, L_w a set of locally updated (or written) TVars, and the stack reflects the depth of the `orElse`-execution; K is a set of TVar-names that is locked by a thread in the commit-phase.

Additional (labeled) variants of the operators `orElse` and `atomically` are required: The operator `orElse!` indicates that `orElse` is active, and the operator `atomically!` indicates that a transaction is active, where `atomically!` has as a second argument an expression that is a saved copy of the start expression and that will again be activated after rollback and restart. I.e., the sets of monadic expressions and MC_{STM} -contexts are adapted as:

$$\begin{aligned}
m \in MExpr \quad & ::= \text{return}_{\text{IO}} e \mid e_1 \gg_{\text{IO}} e_2 \mid \text{future } e \\
& \mid \text{return}_{\text{STM}} e \mid e_1 \gg_{\text{STM}} e_2 \mid \text{atomically } e \\
& \mid \text{atomically! } e\ e' \mid \text{retry} \mid \text{orElse } e_1\ e_2 \mid \text{orElse! } e_1\ e_2 \\
& \mid \text{newTVar } e \mid \text{readTVar } e \mid \text{writeTVar } e
\end{aligned}$$

$$\begin{aligned}
M_{\text{STM}} \in MC_{\text{STM}} \quad & ::= M_{\text{IO}}[\text{atomically! } \widehat{M}_{\text{STM}} e] \\
\widehat{M}_{\text{STM}} \in \widehat{MC}_{\text{STM}} \quad & ::= [\cdot] \mid \widehat{M}_{\text{STM}} \gg_{\text{STM}} e \mid \text{orElse! } \widehat{M}_{\text{STM}} e
\end{aligned}$$

A thread that has a transaction log is called *transactional*, otherwise it is called *non-transactional*. A *CSHF*-process that only has non-transactional threads and no $u\mathbf{t}l\ s$ -components is called *non-transactional*; otherwise, it is called *transactional*.

We say a thread is currently *performing a transaction*, if the current evaluation focusses on the arguments of `atomically!` (see the reduction $\xrightarrow{\text{CSHF}}$ below).

Definition 3.2. *A CSHF-process P is well-formed iff the following holds: Variable names of TVars, threads, and binders are introduced at most once; i.e. threads, bindings, and global TVars are unique per variable. For every component $x\mathbf{t}l\ e$ in a stack-entry of $u\mathbf{t}l\ s$, either there is also a global TVar $x\mathbf{t}g\ e\ o\ g$,*

Monadic IO Computations:	
(lunit _{IO})	$\langle u\!y \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{return}_{\text{IO}} e_1 \gg_{=\text{IO}} e_2] \xrightarrow{\text{CSHF}} \langle u\!y \rangle \Leftarrow \mathbb{M}_{\text{IO}}[e_2 e_1]$
(fork)	$\langle u\!y \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{future } e] \xrightarrow{\text{CSHF}} \nu z, u'. (\langle u\!y \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{return } z] \mid \langle u'\!z \rangle \Leftarrow e)$ where z, u' are fresh and the created thread is not the main thread
(unIO)	$\langle u\!y \rangle \Leftarrow \text{return } e \xrightarrow{\text{SHF}} y = e$ if the thread is not the main-thread
Monadic STM Computations:	
(atomic)	$\langle u\!y \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{atomically } e] \xrightarrow{\text{CSHF}} \nu z. (\langle u\!y \rangle \xleftarrow{\emptyset, [\emptyset]} \mathbb{M}_{\text{IO}}[\text{atomically! } z z] \mid u \text{ tls } [\emptyset] \mid z = e)$
(lunit _{STM})	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{return}_{\text{STM}} e_1 \gg_{=\text{STM}} e_2] \xrightarrow{\text{CSHF}} \langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[e_2 e_1]$
(readl)	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{readTVar } x \mid u \text{ tls } (\{x \text{ tl } e_1\} \dot{\cup} r) : s]$ $\xrightarrow{\text{CSHF}} \nu z. (\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{return}_{\text{STM}} z] \mid z = e_1 \mid u \text{ tls } (\{x \text{ tl } z\} \dot{\cup} r) : s)$
(readg)	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{readTVar } x \mid x \text{ tg } e_1 - g \mid u \text{ tls } r : s]$ $\xrightarrow{\text{CSHF}} \nu z. (\langle u\!y \rangle \xleftarrow{T', L'} \mathbb{M}_{\text{STM}}[\text{return}_{\text{STM}} z] \mid z = e_1 \mid u \text{ tls } (\{x \text{ tl } z\} \dot{\cup} r) : s \mid x \text{ tg } z - g')$ if $x \notin L_a$ where $L = (L_a, L_n, L_w) : L_r, L' = (L_a \cup \{x\}, L_n, L_w) : L_r, T' = T \cup \{x\}$ and $g' = (\{u\} \cup g)$
(writel)	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{writeTVar } x e_1 \mid u \text{ tls } (\{x \text{ tl } e_2\} \dot{\cup} r) : s] \xrightarrow{\text{CSHF}} \langle u\!y \rangle \xleftarrow{T, L'} \mathbb{M}_{\text{STM}}[\text{return}_{\text{STM}} ()] \mid u \text{ tls } (\{x \text{ tl } e_1\} \dot{\cup} r) : s]$ where $L = (L_a, L_n, L_w) : L_r, L' = (L_a, L_n, L_w \cup \{x\}) : L_r$
(writeg)	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{writeTVar } x e_1 \mid x \text{ tg } e_2 - g \mid u \text{ tls } (r : s)]$ $\xrightarrow{\text{CSHF}} \langle u\!y \rangle \xleftarrow{T, L'} \mathbb{M}_{\text{STM}}[\text{return}_{\text{STM}} ()] \mid x \text{ tg } e_2 - g \mid u \text{ tls } (\{x \text{ tl } e_1\} \dot{\cup} r) : s]$ if $x \notin L_a$, where $L = (L_a, L_n, L_w) : L_r, L' = (L_a \cup \{x\}, L_n, L_w \cup \{x\}) : L_r$
(nvar)	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{newTVar } e] \mid u \text{ tls } (r : s) \xrightarrow{\text{CSHF}} \nu x. (\langle u\!y \rangle \xleftarrow{T, L'} \mathbb{M}_{\text{STM}}[\text{return}_{\text{STM}} x] \mid u \text{ tls } (\{x \text{ tl } e\}) \dot{\cup} r : s)$ where $L = (L_a, L_n, L_w)$ and $L' = (\{x\} \cup L_a, \{x\} \cup L_n, L_w)$
(retryup)	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{retry} \gg_{=\text{STM}} e_1] \xrightarrow{\text{CSHF}} \langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{retry}]$
(orElse)	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{orElse } e_1 e_2] \mid (u \text{ tls } (\{x_1 \text{ tl } e_{1,1}, \dots, x_n \text{ tl } e_{1,n}\} : s))$ $\xrightarrow{\text{CSHF}} \nu z_1, \dots, z_n. (\langle u\!y \rangle \xleftarrow{T, L'} \mathbb{M}_{\text{STM}}[\text{orElse! } e_1 e_2] \mid u \text{ tls } ((\{x_1 \text{ tl } z_1, \dots, x_n \text{ tl } z_n\} : (\{x_1 \text{ tl } z_1, \dots, x_n \text{ tl } z_n\}) : s))$ $\mid z_1 = e_{1,1} \mid \dots \mid z_n = e_{1,n})$ where $L = (L_a, L_n, L_w) : L_r, L' = (L_a, L_n, L_w) : ((L_a, L_n, L_w) : L_r)$
(orRetry)	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{orElse! } \text{retry } e_2] \mid u \text{ tls } (r : s) \xrightarrow{\text{CSHF}} \langle u\!y \rangle \xleftarrow{T, L'} \mathbb{M}_{\text{STM}}[e_2] \mid u \text{ tls } s$ where $L = L_e : L'$
(orReturn)	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[\text{orElse! } (\text{return } e_1) e_2] \mid u \text{ tls } (r : s) \xrightarrow{\text{CSHF}} \langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{STM}}[e_1] \mid u \text{ tls } (r : s)$
(retryCGlob)	$\langle u\!y \rangle \xleftarrow{T, L} \mathbb{M}_{\text{IO}}[\text{atomically! } \text{retry } e] \mid x \text{ tg } e_1 - g \xrightarrow{\text{CSHF}} \langle u\!y \rangle \xleftarrow{T', L} \mathbb{M}_{\text{IO}}[\text{atomically! } \text{retry } e] \mid x \text{ tg } e_1 - g'$ if $x \in T \neq \emptyset$, where $T' = T \setminus \{x\}, g' = g \setminus \{u\}$
(retryEnd)	$\langle u\!y \rangle \xleftarrow{\emptyset, L} \mathbb{M}_{\text{IO}}[\text{atomically! } \text{retry } e] \mid u \text{ tls } (r : s) \xrightarrow{\text{CSHF}} \langle u\!y \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{atomically } e]$

Fig. 3. Concurrent Implementation of SHF, transaction reductions

or the TVar is in the L_n -component of the thread-memory (the locally generated TVars). Moreover, for every thread identifier u , there is at most one process component $u \text{ tls } s$. In every stack entry, names of TVars occur at most once. For every thread identifier u in $u \text{ tls } s$ there exists a thread with this identifier.

Though the syntax of *CSHF* is slightly different from *SHF*, we use the same names for the context classes $PC, EC, MC_{\text{STM}}, MC, FC, MC_{\text{IO}}$, and LC_Q . If necessary, then we distinguish the context classes using an index C (for concurrent). For the PC -contexts we assume that also transactional threads are permitted.

Definition 3.3 (Operational Semantics of CSHF). A well-formed *CSHF*-process P reduces to another *CSHF*-process P' (denoted by $P \xrightarrow{\text{CSHF}} P'$) iff $P \equiv D[P_1]$ and $P' \equiv D[P'_1]$ and $P_1 \rightarrow P'_1$ by a reduction rule in Fig. 3, 4, and 5.

We explain the execution of a transaction and the use of the transaction log, where we also point to the rules of the operational semantics.

Start of a transaction (atomic): When the execution is started a new empty transaction log is created.

Read-operations: (readl) and (readg): A read first looks into the local store. If no local TVar exists, then the global value is copied into the local store and the own thread identifier is added to the notify-list of the global TVar.

The commit-phase uses thread-local memory K , written over the thread-arrow after a semicolon.
The third memory-component is the set of locked TVars.

(writeStart) start of commit: locking the read and to-be-written TVars

$$\langle u|y \rangle \xleftarrow{T,L} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x_1 \mathbf{tg} e'_1 - g_1 | \dots | x_n \mathbf{tg} e'_n - g_n$$

$$\xrightarrow{\text{CSHF}} \langle u|y \rangle \xleftarrow{T,L;K} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x_1 \mathbf{tg} e'_1 u g_1 | \dots | x_n \mathbf{tg} e'_n u g_n$$

where $K := \{x_1, \dots, x_n\} = T \cup (L_a \setminus L_n)$ and $L = (L_a, L_n, L_w) : L_r$

(writeClear) removing the notify-entries of u :

$$\langle u|y \rangle \xleftarrow{T,L;K} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x \mathbf{tg} e_2 u g \xrightarrow{\text{CSHF}} \langle u|y \rangle \xleftarrow{T',L;K} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x \mathbf{tg} e_2 u g'$$

if $x \in T$ where $g' = g \setminus \{u\}$ and $T' = T \setminus \{x\}$

(sendRetry) Sending other threads (that are in transactions) a retry:

$$\langle u|y \rangle \xleftarrow{\emptyset,L;K} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x \mathbf{tg} e_2 u g | \langle u'|z \rangle \xleftarrow{T',L'} \mathbb{M}'_{\text{IO}}[\text{atomically!} e_3 e_4]$$

$$\xrightarrow{\text{CSHF}} \langle u|y \rangle \xleftarrow{\emptyset,L;K} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x \mathbf{tg} e_2 u g' | \langle u'|z \rangle \xleftarrow{T',L'} \mathbb{M}'_{\text{IO}}[\text{atomically!} \text{retry } e_4]$$

if $x \in L_w$, $g \neq \emptyset$, $u' \in g$, where $L = (L_a, L_n, L_w) : L_r$, and $g = g' \dot{\cup} \{u'\}$

(writeTV) Overwriting the global TVars with the local TVars of u :

$$\langle u|y \rangle \xleftarrow{\emptyset,L;K} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x \mathbf{tg} e_2 u [] | u \mathbf{tls}(\{x \mathbf{tl} e_3\} \dot{\cup} r : s)$$

$$\xrightarrow{\text{CSHF}} \langle u|y \rangle \xleftarrow{\emptyset,L';K} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x \mathbf{tg} e_3 u [] | u \mathbf{tls}(r : s)$$

if for all $z \in L_w \setminus L_n$ the g is empty in the component $z \mathbf{tg} e_2 u g$, and if $L_w \setminus L_n \neq \emptyset$
where $L = (L_a, L_n, L_w) : L_r$, $x \in L_w \setminus L_n$, $L' = (L_a, L_n, L_w \setminus \{x\}) : L_r$.

(unlockTV) Unlocking the locked TVars:

$$\langle u|y \rangle \xleftarrow{\emptyset,L;K} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x \mathbf{tg} e_2 u [] \xrightarrow{\text{CSHF}} \langle u|y \rangle \xleftarrow{\emptyset,L;K'} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x \mathbf{tg} e_3 - []$$

if $L_w \setminus L_n = \emptyset$, and $K \neq \emptyset$ where $L = (L_a, L_n, L_w) : L_r$, and $K' = K \setminus \{x\}$

(writeTVn) Moving the freshly generated TVars of u into global store:

$$\langle u|y \rangle \xleftarrow{\emptyset,L;\emptyset} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | u \mathbf{tls}(\{x \mathbf{tl} e_3\} \dot{\cup} r : s)$$

$$\xrightarrow{\text{CSHF}} \langle u|y \rangle \xleftarrow{\emptyset,L;\emptyset} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | x \mathbf{tg} e_3 - []$$

if $L_n \neq \emptyset$, $x \in L_n$, where $L = (L_a, L_n, L_w) : L_r$ and $L' = (L_a, L_n \setminus \{x\}, \emptyset) : L_r$

(writeEnd) Removing local store and end of commit of transaction:

$$\langle u|y \rangle \xleftarrow{\emptyset,L;\emptyset} \mathbb{M}_{\text{IO}}[\text{atomically!}(\text{return}_{\text{STM}} e_1) e] | u \mathbf{tls}(r : s) \xrightarrow{\text{CSHF}} \langle u|y \rangle \Leftarrow \mathbb{M}_{\text{IO}}[\text{return}_{\text{IO}} e_1]$$

if no other rules of the commit-phase for u are applicable, i.e., if $L_n = L_w = \emptyset$, where $L = (L_a, L_n, L_w) : L_r$.

Fig. 4. Concurrent Implementation of *SHF*, commit phase of transaction

Functional Evaluation

$$(\text{feval}_{\text{IO}}) P \xrightarrow{\text{CSHF}} P', \text{ if } P \xrightarrow{a} P' \text{ for } a \in \{\text{cp}_{\text{IO}}, \text{abs}_{\text{IO}}, \text{mkb}_{\text{IO}}, \text{lbeta}_{\text{IO}}, \text{case}_{\text{IO}}, \text{seq}_{\text{IO}}\}$$

$$(\text{feval}_{\text{STM}}) P \xrightarrow{\text{CSHF}} P', \text{ if } P \xrightarrow{a} P' \text{ for } a \in \{\text{cp}_{\text{STM}}, \text{abs}_{\text{STM}}, \text{mkb}_{\text{STM}}, \text{lbeta}_{\text{STM}}, \text{case}_{\text{STM}}, \text{seq}_{\text{STM}}\}$$

The reductions with parameter $Q \in \{\text{STM}, \text{IO}\}$ are defined as follows

$$(\text{cp}_Q) \quad \mathbb{L}_Q[x_1 \mid x_1 = x_2 \mid \dots \mid x_{n-1} = x_n \mid x_n = v] \rightarrow \mathbb{L}_Q[v \mid x_1 = x_2 \mid \dots \mid x_{n-1} = x_n \mid x_n = v],$$

if v is an abstraction, a cx-value, or a component name

$$(\text{abs}_Q) \quad \mathbb{L}_Q[x_1 \mid x_1 = x_2 \mid \dots \mid x_{m-1} = x_m \mid x_m = c \ e_1 \dots e_n]$$

$$\rightarrow \nu y_1, \dots, y_n. (\mathbb{L}_Q[c \ y_1 \dots y_n \mid x_1 = x_2 \mid \dots \mid x_{m-1} = x_m \mid x_m = c \ y_1 \dots y_n \mid y_1 = e_1 \mid \dots \mid y_n = e_n])$$

if c is a constructor, or `returnSTM`, `returnIO`, `>>=STM`, `>>=IO`, `orElse`, `atomically`, `readTVar`, `writeTVar`, `newTVar`, or `future`, and $n \geq 1$, and some e_i is not a variable.

$$(\text{mkb}_Q) \quad \mathbb{L}_Q[\text{letrec } x_1 = e_1, \dots, x_n = e_n \text{ in } e] \rightarrow \nu x_1, \dots, x_n. (\mathbb{L}_Q[e \mid x_1 = e_1 \mid \dots \mid x_n = e_n])$$

$$(\text{lbeta}_Q) \quad \mathbb{L}_Q[(\lambda x. e_1) \ e_2] \rightarrow \nu x. (\mathbb{L}_Q[e_1 \mid x = e_2])$$

$$(\text{case}_Q) \quad \mathbb{L}_Q[\text{case}_T (c \ e_1 \dots e_n) \text{ of } \dots ((c \ y_1 \dots y_n) \rightarrow e) \dots] \rightarrow \nu y_1, \dots, y_n. (\mathbb{L}_Q[e \mid y_1 = e_1 \mid \dots \mid y_n = e_n]), \text{ if } n > 0$$

$$(\text{case}_Q) \quad \mathbb{L}_Q[\text{case}_T c \text{ of } \dots (c \rightarrow e) \dots] \rightarrow \mathbb{L}_Q[e]$$

$$(\text{seq}_Q) \quad \mathbb{L}_Q[(\text{seq } v \ e)] \rightarrow \mathbb{L}_Q[e], \text{ if } v \text{ is a functional value}$$

Fig. 5. Concurrent implementation of *SHF*, functional reductions

Write-operations (writel) and (writeg): A write command always writes into the local store, perhaps preceded by a copy from the global into the local store .

OrElse-Evaluation: If evaluation descends into the left expression of `orElse`, then the stacks of TVars and of local TVars are extended by duplicating their top element in rule (`orElse`). For the final write in the commit phase only the top of the stack is relevant. If evaluation of the left expression is successful, the stack remains as it is (by rule (`orReturn`)). In case of a retry, the top element is popped (by rule (`orRetry`)) and the execution of the second expression then uses the stack before executing the `orElse`. Note that the information on the read TVars is kept in the set T , since the values may have an influence on the outcome of the `orElse`-tree execution. This is necessary, since the big-step semantics of *SHF* enforces a left-to-right evaluation of the `orElse`-tree, where the leftmost successful try will be kept. Note that this is semantically different from making a non-deterministic choice of one of the successful possibilities in the `orElse`-tree.

Commit-Phase: At the end of a transaction, there is a lock-protected sequence of updates: A thread that is in its commit-phase first locks all its globally read and to-be-updated TVars (rule (`writeStart`)). Locked variables cannot be accessed by other threads for reading, or modifying the notify-list. Then all own notification-entries are removed (rule (`writeClear`)). All the threads, which are unequal to the running thread, and that are in the notify-list of the updated TVars, will be stopped by replacing the current transaction expression by `retry` (rule (`sendRetry`)). This mechanism is like raising (synchronous) exceptions to influence other threads. The to-be-updated TVars are written into the global store (rule (`writeTV`)), then the locks are released (rule (`unlockTV`)) and fresh TVars are also moved to the global store (rule (`writeTVn`)). Finally, the transaction log is removed (rule (`writeEnd`)).

Rollback and Restart: A transaction is rolled back and restarted by the `retry`-command (if it is not inside an `orElse`-command). This can occur by a user programmed `retry`, or if the transaction gets stopped by a conflicting transaction which is committing. The thread removes the notification entries (rule (`retryCGlob`)) and then the transaction code is replaced by the original expression (rule (`retryEnd`)).

Also, it is easy to extend the monomorphic type system to *CSHF* and to see that reduction keeps well-typedness.

Definition 3.4. A *CSHF*-process P is successful, iff it is well-formed and the main thread is of the form $\langle u!y \rangle \xrightarrow{\text{main}} \text{return } e$. May- and should-convergence in *CSHF* are defined by: $P \downarrow_{\text{CSHF}}$ iff P is well-formed and $\exists P' : P \xrightarrow{\text{CSHF},*} P' \wedge P'$ successful. $P \Downarrow_{\text{CSHF}}$, iff P is well-formed and $\forall P' : P \xrightarrow{\text{CSHF},*} P' \implies P' \downarrow_{\text{CSHF}}$. Must- and may-divergence of process P are the negations of may- and should-convergence and are denoted by $P \uparrow_{\text{CSHF}}$, and $P \Uparrow_{\text{CSHF}}$, resp., where $P \uparrow_{\text{CSHF}}$ is also equivalent to $P \xrightarrow{\text{CSHF},*} P'$ such that $P' \uparrow_{\text{CSHF}}$.

Contextual approximation \leq_{CSHF} and equivalence \sim_{CSHF} in *CSHF* are defined as $\leq_{\text{CSHF}} := \leq_{\downarrow_{\text{CSHF}}} \cap \leq_{\Downarrow_{\text{CSHF}}}$ and $\sim_{\text{CSHF}} := \leq_{\text{CSHF}} \cap \geq_{\text{CSHF}}$ where for $\zeta \in \{\downarrow_{\text{CSHF}}, \Downarrow_{\text{CSHF}}\}$: $P_1 \leq_{\zeta} P_2$ iff $\forall \mathbb{D} \in \text{PC}_C : \mathbb{D}[P_1]\zeta \implies \mathbb{D}[P_2]\zeta$.

4 Correctness of the Concurrent Implementation

In this section we show that *CSHF* can be used as a correct evaluator for *SHF* and its semantics. Hence, we provide a translation from *SHF* into *CSHF*:

Definition 4.1. *The translation ψ of an SHF-process into a CSHF-process is defined homomorphically on the structure of processes: Usually it is the identity on the constructs; the only exception is $\psi(x \mathbf{t} e) := x \mathbf{t} g e - []$, i.e. initially, the list of threads to be notified is empty and the TVar is not locked. CSHF-processes $\psi(P)$ where P is an SHF-process are the initial CSHF-processes.*

We are mainly interested in *CSHF*-reductions that start with initial *CSHF*-processes. Since only transactions can introduce local TVars which are removed at the end of a transaction, the following lemma holds:

Lemma 4.2. *Every initial CSHF-process is well-formed provided the corresponding SHF-process is well-formed. Also, every reduction descendant of an initial CSHF-process is well-formed.*

The correctness theorem we want to prove is the following:

Main Theorem 4.3. *The translation $\psi : SHF \rightarrow CSHF$ is observational correct, i.e. for all process contexts D and SHF-processes P the equivalences $D[P] \Downarrow \iff \psi(D)(\psi(P)) \Downarrow_{CSHF}$ and $D[P] \Downarrow \iff \psi(D)(\psi(P)) \Downarrow_{CSHF}$ hold. This also implies that ψ is adequate, i.e. $\psi(P_1) \sim_{CSHF} \psi(P_2) \implies P_1 \sim_c P_2$.*

Adequacy is a direct consequence of observational correctness (see [SSNSS08]). Since the translation ψ is compositional, i.e. $\psi(D)(\psi(P)) = \psi(D[P])$ for any processes context D and process P , for the proof of observational correctness it is sufficient to show that ψ preserves and reflects may- and should convergence:

Theorem 4.4. *The translation $\psi : SHF \rightarrow CSHF$ is convergence equivalent, i.e. for all SHF-processes P the equivalences $P \Downarrow \iff \psi(P) \Downarrow_{CSHF}$ and $P \Downarrow \iff \psi(P) \Downarrow_{CSHF}$ hold.*

In the remainder of this section we will prove this theorem to complete the proof of Main Theorem 4.3. Thus we have to show that may- and should-convergence are the same for the big-step semantics and the concurrent implementation.

In order to be on solid ground, we first analyze the invariants during transactions and properties of the valid configurations.

Lemma 4.5. *The following properties hold during a CSHF-reduction on a well-formed CSHF-process that is reachable from a non-transactional CSHF-process.*

1. *For every component $x \mathbf{t} e$, either $x \in L_n$ of the top entry in L , or there is a global TVar x . For every pair of thread identifier u , and TVar-name x , every stack element of the TVar-stack for u contains at most one entry $x \mathbf{t} e$.*
2. *If u is in a notify-list of a global TVar, then thread u is transactional.*
3. *Every transaction that starts the commit-phase for thread u by performing the rule (*writeStart*) is able to perform all other rules until (*writeEnd*) is performed, without retry, nontermination or getting stuck.*

We require the notion of convergence equivalence for program transformations, which is analogously defined to convergence equivalence of translations:

Definition 4.6. *If $P_1 \Downarrow_{CSHF} \iff P_2 \Downarrow_{CSHF}$ and $P_1 \Downarrow_{CSHF} \iff P_2 \Downarrow_{CSHF}$ then we write $P_1 \sim_{ce} P_2$. A program transformation ξ (i.e. a binary relation over CSHF-processes) is convergence equivalent iff $P_1 \xi P_2$ always implies $P_1 \sim_{ce} P_2$.*

4.1 Convergence Equivalence of Special Transformations

In this section we show convergence-equivalence of special transformations in the concurrent calculus *CSHF*. These results will be used in the preceding subsection to proof correctness of the translation.

Since for the rules of the calculus *CSHF*, the variable (binding-) chains are transparent, there is no need to copy binding-variables to other places in the expressions, which strongly reduces the number of critical overlappings. This, however, enforces that (abs) is in the set of standard reductions.



Fig. 6. The standard: square and triangle diagrams of special transformations in Fig. 9

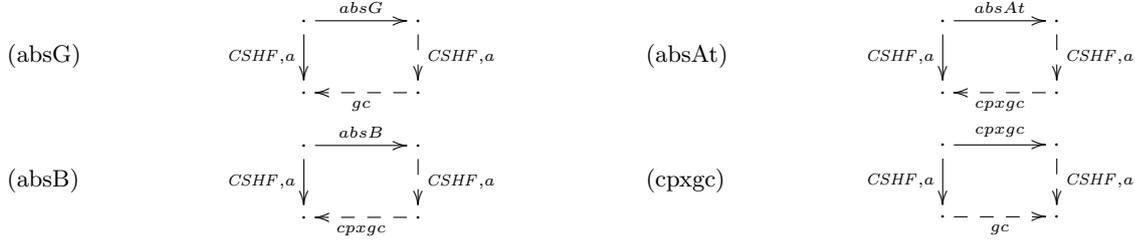


Fig. 7. The unusual forking diagrams of special transformations in Fig. 9

In the following we use forking and commuting diagrams, where in general we only write the forking diagrams. A forking diagram shows (in an abstract way, since processes are omitted) how a given overlapping $\leftarrow \xrightarrow{CSHF} \xrightarrow{T}$ (the fork) between a standard reduction and a transformation T can be closed, a commuting diagram shows how a sequence $\xrightarrow{T} \xrightarrow{SHF}$ can be closed, i.e. how a transformation and a reduction can be commuted. A set of diagrams for a transformation T is complete if for every concrete fork $P_1 \xleftarrow{CSHF} P_2 \xrightarrow{T} P_3$ (or sequence $P_1 \xrightarrow{T} P_2 \xrightarrow{SHF} P_3$) at least one diagram of the set is applicable, which means that the corresponding transformations and reductions exist (on the concrete level). Since there may be several concurrent threads, the diagrams have to express reduction commutations of concurrent reductions. In order to have simple forms of diagrams, at most a single sr-reduction should be vertically in the diagrams on the west- and east-side. If more standard reductions would be necessary, then reasoning is far more complex.

Forking with a (atomic) standard reduction:

$$\begin{array}{ccc}
 \langle u|x \rangle \leftarrow \text{atomically } e & \xrightarrow{\text{absAt}} & \langle u|x \rangle \leftarrow \text{atomically } z \\
 \downarrow \text{CSHF, (atomic)} & & \downarrow \text{CSHF, (atomic)} \\
 \langle u|x \rangle \leftarrow \text{atomically! } z \mid z = e & \xleftarrow{\text{cpxgc}} & \langle u|x \rangle \leftarrow \text{atomically! } z' \mid z' = z \mid z = e
 \end{array}$$

Forking standard reduction (cp) with (cpxgc):

$$\begin{array}{ccc}
 \nu x.E[x] \mid x = y \mid y = v & \xrightarrow{\text{cpxgc}} & \nu x.E[y] \mid y = v \\
 \downarrow \text{CSHF, cp} & & \downarrow \text{CSHF, cp} \\
 \nu x.E[v] \mid x = y \mid y = v & \xleftarrow{\text{gc}} & \nu x.E[v] \mid y = v
 \end{array}$$

Forking (writeTV) with (absG)

$$\begin{array}{ccc}
 x \text{ tg } e \text{ o g } \dots & \xrightarrow{\text{absG}} & x \text{ tg } z \text{ o g } \mid z = e \\
 \downarrow \text{CSHF, writeTV} & & \downarrow \text{CSHF, writeTV} \\
 x \text{ tg } e' \text{ o g } \dots & \xleftarrow{\text{gc}} & x \text{ tg } e' \text{ o g } \mid z = e
 \end{array}$$

Forking (cp) with (absB)

$$\begin{array}{ccc}
 E[x] \mid x = c \ y_1 \ y_2 \dots & \xrightarrow{\text{absB}} & E[x] \mid x = c \ z_1 \ z_2 \mid z_1 = y_1 \mid z_2 = y_2 \\
 \downarrow \text{CSHF, cp} & & \downarrow \text{CSHF, cp} \\
 E[c \ y_1 \ y_2] \mid x = c \ y_1 \ y_2 \dots & \xleftarrow{\text{cpxgc}} & E[c \ z_1 \ z_2] \mid x = c \ z_1 \ z_2 \mid z_1 = y_1 \mid z_2 = y_2
 \end{array}$$

Fig. 8. Examples for diagrams in Lemma 4.7

Lemma 4.7. *The forking diagrams of the transformations in Fig. 9 with standard-reductions may be square or triangle diagrams or the equality (see Fig. 6). The unusual forking diagrams of the transformations in Fig. 9 are in Fig. 7. The commuting diagrams can be obtained from them by taking the same diagram where the existential and given vertical arrows are switched. The equality diagram may only occur as follows: an (absB)-transformation may be a (CSHF,abs)-reduction, a (cpBE)-transformation may be a (CSHF,cp)-reduction, and a (funrB)-transformation may be a functional standard reduction.*

Proof. We provide arguments for every reduction, and also exhibit exceptional diagrams in Fig. 8. For a (gc) transformation it is easy to verify that bindings cannot be removed by any standard reduction, and standard reductions do not interfere with the removed bindings by (gc). For (cpBE) the restriction of the target positions shows that there is no fork, where the to-be-copied expression of a standard reduction is modified by (cpBE). The sharing mechanism in, for example, (orElse) and (readg) shows that there is no duplicate of (cpBE) in the south edge of diagrams. A further argument is that the standard reduction is deterministic within threads. For (funrB) the reduction takes place inside a binding and never inside a thread, hence standard reductions and (funrB) only have trivial overlappings. The (absG) transformation may critically overlap with a (writeTV) standard reduction. In this case a (gc) is required to remove the created bindings. The (absAt) transformation can critically overlap with the (atomic) reduction where (cpxgc) is necessary to remove indirections. The transformation (absB) may critically overlap with a (cp) that copies a cx-value where again (cpxgc) can be used to remove indirections. The exceptional diagram for (cpxgc) using (gc) in the south arrow occurs for example in a forking of (case), (seq), or (lbeta) with (sendRetry).

Lemma 4.8. *Given a reduction sequence of the form $P_1(\xrightarrow{CSHF} \cup \xrightarrow{spec})^* P_2$, where \xrightarrow{spec} is a union of (cpBE), (absB), (funrB), (absG), (absAt), (gc), (cpxgc) and their inverses, then the reduction can be rearranged as $P_1 \xrightarrow{CSHF,*} P_2 \xrightarrow{spec,*} P_2$.*

Proof. This follows from the forking diagrams for the special transformations and for (cpxgc), as well as the commuting diagrams, which can be obtained from the forking diagrams.

Lemma 4.9. *The rules in Fig. 9 are convergence equivalent for may-convergence and should-convergence in the calculus CSHF.*

Proof. We first consider may-convergence. Note that due to the reverse transformation steps in the diagrams, it is necessary to show that the transformations preserve and reflect convergence in one induction. Let $P' \xleftarrow{CSHF,k} P \xrightarrow{\tau} P''$ where τ is any of the special transformations or their inverse transformation, P' is successful and $k \geq 0$. By induction on k we show that P'' is may-convergent. For the base case, P is already successful. Then P'' must be successful, too, which follows by inspection of the definitions of the transformations. For the induction step assume $P' \xleftarrow{CSHF,k} P_1 \xleftarrow{CSHF} P \xrightarrow{\tau} P''$ for some $k \geq 0$. Then depending on the direction of τ we apply a forking or a commuting diagram of Lemma 4.7 to $P_1 \xleftarrow{CSHF} P \xrightarrow{\tau} P''$. Then there are the following cases:

- If the second or the third diagram of Fig. 6 is applied, then either $P'' \xrightarrow{CSHF} P_1$ or $P'' = P_1$. In both cases this implies $P'' \downarrow_{CSHF}$.
- If any other diagram is applied, then there exists a process P_2 and a special transformation or its inverse τ' such that $P_1 \xrightarrow{\tau'} P'' \xrightarrow{CSHF} P_2$. Now we apply the induction hypothesis to $P' \xleftarrow{CSHF,k} P_1 \xrightarrow{\tau'} P_2$ and derive $P_2 \downarrow_{CSHF}$. Since $P'' \xrightarrow{CSHF} P_2$, we also have $P'' \downarrow_{CSHF}$.

The proof for should-convergence is completely analogous except for the base case: The base case requires that if $P \xrightarrow{\tau} P'$ then $P \uparrow_{CSHF} \iff P' \uparrow_{CSHF}$. But this holds by equivalence w.r.t. may-convergence.

For proving convergence equivalence of ψ (Theorem 4.4) we have to show four parts: preservation of may-convergence ($P \downarrow \Rightarrow \psi(P) \downarrow_{CSHF}$), reflection of may-convergence ($\psi(P) \downarrow_{CSHF} \Rightarrow P \downarrow$), preservation of should-convergence ($P \downarrow \Rightarrow \psi(P) \downarrow_{CSHF}$), and reflection of should-convergence ($\psi(P) \downarrow_{CSHF} \Rightarrow P \downarrow$).

4.2 Preservation of May-Convergence

We have to show that $P \downarrow$ implies $\psi(P) \downarrow_{CSHF}$, which is not completely straightforward, since the big-step reduction $\xrightarrow{SHFA,*}$ can be tried for free, and only in the case of success, i.e., a `returnSTM` is obtained by an

- (cpBE) $x = \mathbb{E}[x_1] \mid x_1 = x_2 \mid \dots \mid x_{m-1} = x_m \mid x_m = v \longrightarrow x = \mathbb{E}[v] \mid x_1 = x_2 \mid \dots \mid x_{m-1} = x_m \mid x_m = v$,
if v is an abstraction, a cx-value or a component-name, and $\mathbb{E} \neq [\cdot]$.
- (absB) $x = c \ e_1 \dots e_n \longrightarrow \nu x_1, \dots, x_n. x = c \ x_1 \dots x_n \mid x_1 = e_1 \mid \dots \mid x_n = e_n$.
- (funrB) Every functional rule (without the surrounding \mathbb{L} -context) in a context $x = \mathbb{E}[\cdot]$, but not (cp) and not (abs).
- (absG) $x \mathbf{t}g \ e \ o \ g \longrightarrow \nu z. x \mathbf{t}g \ z \ o \ g \mid z = e$ where $o \in \{-, u\}$
- (absAt) $\langle x \lambda u \rangle \Leftarrow \mathbb{M}_{\text{IO}}[(\mathbf{atomically} \ e)] \longrightarrow \langle x \lambda u \rangle \Leftarrow \mathbb{M}_{\text{IO}}[(\mathbf{atomically} \ z)] \mid z = e$.
- (gc) $\nu x_1, \dots, x_n. P \mid x_1 = e_1 \mid \dots \mid x_n = e_n \longrightarrow \nu x_1, \dots, x_n. P$ if for all $i = 1, \dots, n$: x_i does not occur free in P .
- (cpXgc) $\nu x. C[x, \dots, x] \mid x = y \longrightarrow C[y, \dots, y]$
where C is a multi-context, where all holes are in an \mathbb{A} -context, and all occurrences of x are indicated in the notation, i.e., there are no other occurrences. \mathbb{A} is the class of expression contexts where the hole is not within an abstraction nor in a letrec-expression, nor in an alternative of a case and not in an argument of a constructor.

Fig. 9. Special Transformations for *CSHF*

atomic transaction, the changes of the $\xrightarrow{SHFA,*}$ -reduction sequence are kept. Also, if an `orElse`-expression is reduced, the big-step reduction is permitted to evaluate the second expression, if it is known that the first one would end in a retry. This is different in *CSHF*, since the execution has to first evaluate the left expression of an `orElse`-expression, and only in case it “retries”, the second expression will be evaluated where the changes of the TVars are not kept, but changes belonging to functional evaluation in the bindings are kept. Analyzing the behavior exhibits that these changes of the process can be proved as convergence equivalent transformations.

As a base case, the following lemma holds:

Lemma 4.10. *If P is successful, then $\psi(P)$ is successful.*

An easy case are non-(atomic)-standard reductions:

Lemma 4.11. *If $P_1 \xrightarrow{SHF,a} P_2$ where $a \neq (\mathbf{atomic})$, then $\psi(P_1) \xrightarrow{CSHF} \psi(P_2)$.*

The more complex cases arise for the transactions in the big-step reduction. In this case the state of the concurrent implementation consists of stacks, global TVars and stacks of local TVars. We consider several program transformations related to reduction rules, which are chosen such that it is sufficient to simulate the modifications in the bindings of the retried *CSHF*-standard reductions and then to rearrange reduction sequences of the concurrent implementation.

Definition 4.12 (CSHF special transformations). *The special transformations are defined in Fig. 9 where we assume that they are closed w.r.t. \mathbb{D} -contexts and structural congruence, i.e. for any transformation \xrightarrow{a} with $a \in \{(\mathbf{cpBE}), (\mathbf{absB}), (\mathbf{funrB}), (\mathbf{absG}), (\mathbf{absAt}), (\mathbf{gc})\}$ we extend its definition as follows: If $P \equiv D[P']$, $Q \equiv D[Q']$, and $P' \xrightarrow{a} Q'$, then also $P \xrightarrow{a} Q$.*

First we observe the effects of global retries:

Lemma 4.13. *If there is a *CSHF*-reduction sequence $P_1 \xrightarrow{CSHF,*} P_2$, where an STM-transaction is started in the first reduction for thread u , and the last reduction is a global retry, i.e. (`retryEnd`) for thread u , of the transaction, then $P_1 \left(\xrightarrow{CSHF} \cup \xrightarrow{spt} \right)^* P_2$ where \xrightarrow{spt} is the union of the *CSHF*-special-transformations (`cpBE`), (`absB`), (`funrB`), (`absG`), and inverse (`gc`)-transformations from Definition 4.12.*

Proof. A global retry removes all generated local TVars for this thread. There may remain changes in the global TVars and in the bindings: Transformations (`absG`) may be necessary for global TVars. Functional transformations in the sharing part are still there, but may be turned into non-standard reductions, if these were triggered only by the thread u . Since (`cp`)-effects in the thread expression are eliminated after a retry in an `orElse`: These may be (`cpBE`), (`absB`), (`funrB`). Various reductions generate bindings which are no longer used after removal of the first expression in an `orElse`, which can be simulated by a reverse (`gc`). These special reductions replace the transaction reductions for thread u , but the others remain, hence the lemma holds.

Lemma 4.14. *If $P_1 \xrightarrow{SHF,atomic} P_2$, then $\psi(P_1) \xrightarrow{CSHF,*} P'_2$, and $P'_2 \xrightarrow{sptx,*} \psi(P_2)$, where \xrightarrow{sptx} consists of special transformations in Fig. 9 and their inverses.*

Proof. The correspondence between the contents of the global TVar x in $CSHF$ for a single thread u is the local TVar x on the top of the stack, or the global TVar if there is no local TVar for x . The rules that change the bindings permanently in the atomic-transaction reduction are: (atomic), (readl), (readg), (orElse), which can be simulated by sequences of (absG) and reverse (gc). The effects of the functional rules ((cp_Q), (abs_Q), (mkb_Q), (lbeta_Q), (case_Q)) that survive a retry are either simulated by (mkb_Q), (lbeta_Q), (case_Q) in a binding, or by (cpBE), (absB), or inverse (gc). An (ortry)-reduction in SHF can be simulated in $CSHF$ by a sequence of reductions starting with (orElse) and ending with (orRetry). However, since (ortry) undoes all changes, in $CSHF$ it is necessary to undo the changes in bindings by the special transformations.

The rearrangement is possible as claimed using Lemma 4.8.

Theorem 4.15. *For all SHF-processes P , we have $P \downarrow \implies \psi(P) \downarrow_{CSHF}$.*

Proof. This follows by an induction on the length of the given reduction sequence for P . The base case is covered in Lemma 4.10. Lemmas 4.11, 4.14, and 4.9 show that if $P_1 \xrightarrow{SHF} P_2$, then there is some $CSHF$ -process P'_2 such that $\psi(P_1) \xrightarrow{CSHF,*} P'_2$ with $P'_2 \sim_{ce} \psi(P_2)$. This is sufficient for the induction step.

4.3 Reflection of May-Convergence

In this section we distinguish the different reduction steps within a $CSHF$ -reduction sequence $\psi(P_1) \xrightarrow{*} P_2$ for the different threads u . A (sendRetry)-reduction belongs to the sending thread. A subsequence of the reduction sequence starting with (atomic) and ending with (writeEnd), which includes exactly the u -reduction steps in between, and where no other (atomic) or (writeEnd)-reductions are contained is called a *transaction*. A prefix of a transaction is also called a *partial transaction*. The subsequence of an u -transaction for thread u starting with (writeStart) and ending with (writeEnd) is called the *commit-phase* of the transaction. If the subsequence for thread u starts with (atomic) and ends with (retryEnd), without intermediate (atomic) or (retryEnd), then it is an *aborted transaction*. The subsequence from the first (retryCglobe) until (retryEnd) is the *abort-phase* of the aborted transaction. A prefix of an aborted transaction is also called a *partial aborted transaction*.

Theorem 4.16. *For all SHF-processes P , we have $\psi(P) \downarrow_{CSHF} \implies P \downarrow$.*

Proof. Let $\psi(P_1) \xrightarrow{CSHF,*} P_2$ where P_2 is successful. Since the transactions in the reduction sequence may be interleaved with other transactions and reduction steps, we have to rearrange the reduction sequence in order to be able to retranslate it into SHF .

Partial Transactions that do not contain a (writeStart) within the reduction sequence can be eliminated and thereby replaced by interspersed special transformations using Lemma 4.13, which again can be eliminated from the successful reduction sequence by Lemma 4.9. If the partial transaction contains a (writeStart), then the missing reduction steps can be added within the reduction sequence before a successful process is reached, since the commit-phase does not change the successful-property of processes. *Aborted Transactions* can be omitted since they are replaceable by special transformations, which again can be removed.

Grouping Transactions: We can now assume that within the reduction sequence $\psi(P_1) \xrightarrow{CSHF,*} P_2$ where P_2 is successful, all transactions are completed, and that there are no aborted ones. Now we rearrange the reduction sequence: The (writeEnd)-reduction step is assumed to be the point of attraction for every transaction. Moving single reduction steps starts from the rightmost non-grouped transaction. For this u -transaction, we move the reduction steps that belong to it in the direction of its (writeEnd), i.e., to the right. The reduction steps that belong to the same transaction and are between (writeStart) and (writeEnd) can be moved to the right squeezing out the non- u -reduction steps, which is possible, since the locking prevents other transaction reduction steps of other threads to be in between, and since there are no functional reduction steps in between.

Now we speak of the u -block of reduction steps, if every reduction step in it belongs to the same transaction and ends with (writeStart). The block is complete, if the first (leftmost) reduction step is (atomic). The interesting part of the reduction sequence is like $\xrightarrow{a}; S; B^u$, where B^u is the u -block; S is a reduction sequence of non- u -reduction steps and \xrightarrow{a} is the reduction step that we want to move to the right, before B^u . For the reduction steps for u that are before (writeStart), we distinguish functional and non-functional reductions. Every standard functional reduction is triggered (or requested) by one or more threads. This must be taken into account when reasoning about shifting functional transactional reduction

steps within reduction sequences. Another issue is which reductions are in S . Since we started with the rightmost transaction, the steps in S are either non-transactional ones, or if they are transactional ones, then these are for another thread. It is obvious that S does not contain a commit-phase of a transaction that aborts u . Thus \xrightarrow{a} cannot be a u -read of a TVar that is committed in S ; it can be a u -write of a TVar, even if writing this TVar is committed in S , since there is no registration of threads at written TVars, and so the rules can be interchanged. Hence, if \xrightarrow{a} is a non-functional transaction step, then it can be moved to the u -block. If \xrightarrow{a} is a functional transaction that is only requested by u , then it will also be moved to the u -block. Now assume that \xrightarrow{a} is a functional transaction that is also requested by another transaction of thread u' . If all such threads are either u or their blocks are to the right of the u -block, then we move the reduction in front of the u -block. Otherwise, i.e, if at least one requesting transaction is in S , then \xrightarrow{a} is now a part of S .

Thus the general situation is: the u -block is the current focus of shifting and there may be more blocks to the right of this block; S consists of functional reductions requested by threads in S , and of other reduction steps. Also in this general case the following holds: functional reductions \xrightarrow{a} that are only requested by u , can be moved to the u -block. Finally, there will be an (atomic)-reduction, and after the move, we have a complete u -block that corresponds to a complete transaction. If the moving is done exhaustively, then the reduction sequence is almost in a form that can be backtranslated into SHF . Every reduction steps that are not within transactions for a thread can be backtranslated as a SHF -reduction.

The last issue in the backtranslation are the local retries in `orElse`-expressions in $CSHF$, which do not have an effect on the TVars, but may apply functional reductions to the bindings, however, in SHF the `orElse`-retries are without any effect.

We use Lemma 4.9 to correctly backtranslate transactions, in particular (orElse)-reductions, including the retries in the transaction. The backtranslation is done from left to right. The plan is to show that the `orElse`-computation tree that is done sequentially in $CSHF$, will also be performed in the same sequence in SHF . A difference are that $\xrightarrow{SHFA,*}$ -reductions are part of the reductions in $CSHF$. Therefore, consider the case that evaluation of the left expression e_1 in a `orElse` ends in a retry. In $CSHF$, all the functional reductions in the reduction sequence corresponding to e_1 can be shifted in the reduction sequence to the right end of the complete reduction sequence, which can then be ignored leaving the sequence successful. Using induction on the depth of the `orElse`-expressions, we can show that there is a backtranslation of the transaction for thread u as an atomic reduction in SHF .

From a bird eyes point of view, a successful $CSHF$ -reduction of a process P is first rearranged into another a successful $CSHF$ -reduction, and then an interleaved rearrangement and backtranslation leads to a converging SHF -reduction sequence for P as a SHF -process.

Corollary 4.17. *Let P be an SHF -process. Then $P \uparrow \iff \psi(P) \uparrow_{CSHF}$.*

4.4 Reflection of Should-Convergence

Instead of proving reflection of should-convergence directly, we equivalently show preservation of may-divergence. The proof is analogous to the preservation of may-convergence in Theorem 4.15, where now, however, the reduction sequence ends in a must-divergent process. Using Corollary 4.17 as a base case shows that may-divergence is preserved:

Theorem 4.18. *For all SHF -processes P , we have $P \uparrow \implies \psi(P) \uparrow_{CSHF}$.*

4.5 Preservation of Should-Convergence

This is the last part of the analysis. A needed lemma is the following which shows that partial transactions can be eliminated:

Lemma 4.19. *Let P_0 be a nontransactional process, $P_0 \xrightarrow{CSHF,*} P_1$ with $P_1 \uparrow_{CSHF}$, such that $P \xrightarrow{Q} P_1$ is a suffix of the reduction and $P \xrightarrow{Q} P_1$ is a partial transaction for a thread u starting with (atomic) but the commit-phase and the retry-phase is missing, i.e. there is no reduction (`writeStart`) nor a (`retryCGlob`) for thread u . Then $P \uparrow_{CSHF}$.*

Proof. The proof is by contradiction. Assume that $P \downarrow_{CSHF}$, i.e. $P \xrightarrow{CSHF,*} P_2$, where P_2 is successful. Then we have to argue that this implies $P_1 \downarrow_{CSHF}$.

$$\begin{array}{ccc}
 P & \xrightarrow[CSHF]{Q,u} & P_1 \uparrow \\
 \downarrow & & \downarrow \\
 P_2 \text{ (success)} & & \cdot
 \end{array}$$

The assumption that P_0 is a nontransactional process implies that the reduction sequences are reachable and thus the notifications and memory at the thread arrows and the TVars are consistent.

Using the claims in the proof of Theorem 4.16, we can assume that in the reduction sequence $P \xrightarrow{CSHF,*} P_2$ all transactions are completed and that there are no aborted transactions.

The goal is now to construct a converging *CSHF*-reduction sequence for P_1 . The idea is to use the same reduction steps as for P . Let us ignore the functional and the *IO*-reduction steps and concentrate on the transactional ones. The reduction steps can be shifted from $P \xrightarrow{CSHF,*} P_2$ to a reduction sequence for P_1 . Note that this shifting may also be accompanied by modified notification entries of TVars. There are two cases:

1. If the TVars read by the transaction Q are unchanged in $P \xrightarrow{CSHF,*} P_2$ (up to bindings via variable-chains), then a construction of a reduction of P_1 to a successful process is possible in a standard way, where the u -transaction is either unused or will be completed.
2. Let V be the set of global TVars read by the transaction Q . Assume some TVar from V is updated in the prefix of $P \xrightarrow{CSHF,*} P_2$, before any u -transaction starts. then we can again construct a converging standard reduction for P_1 : Transporting the reduction steps to P_1 , the only difference are the notifications for the TVars read by the u -transaction. There will be an extension by a (*sendRetry*) that aborts the (image of the) Q -transaction, which we will add immediately after a commit of the updates.

After the retry the state (i.e. the process) is the same in the two reductions up to some special transformations. Now we can proceed with the construction that transports the reductions from the successful reduction sequence to the P_1 -reduction sequence. This will end in a converging standard reduction sequence for P_1 .

This contradicts the assumption that P_1 is must-divergent.

Now we show reflection of of may-divergence, which is equivalent to preservation of should-convergence.

Theorem 4.20. *For all SHF-processes P , we have $\psi(P) \uparrow_{CSHF} \implies P \uparrow$.*

Proof. Assume given a non-transactional *CSHF*-process $\psi(P)$ with $\psi(P) \xrightarrow{CSHF,*} P_1$ and $P_1 \uparrow_{CSHF}$. The reasoning is as in Theorem 4.16 with some differences, since we have to ensure the condition $P_1 \uparrow_{CSHF}$, which is more complex than the “successful”-condition. *Partial transactions:* If it includes a (*writeStart*) or a (*retryCGlob*), then the transaction can be completed by $P_1 \xrightarrow{CSHF,*} P'_1$ where $P'_1 \uparrow_{CSHF}$, and so we can assume that these do not exist. If the partial transaction does neither include a (*writeStart*) nor a (*retryCGlob*), then using the same arguments as in Theorem 4.16, we see that we can assume that the partial transaction is grouped before P_1 , i.e. the transaction ends with the partial transaction $P''_1 \xrightarrow{CSHF,*} P_1$. Lemma 4.19 shows that $P''_1 \uparrow_{CSHF}$, which permits to assume that the partial transaction can be omitted. *Aborted transactions:* can be omitted since they are replaceable by special transformations, which again can be removed due to Lemma 4.9. *Grouping transaction:* same arguments as in the proof of Theorem 4.16. *Final retranslation:* similar to the proof of Theorem 4.16.

4.6 Summary

We have proved in Theorems 4.15, 4.16, 4.18, and 4.20 that the translation ψ mapping *SHF*-processes into *CSHF*-processes is convergence equivalent. Thus we have proved Theorem 4.4 and thus also Main Theorem 4.3 (observational correctness and adequacy of ψ), since ψ is compositional. This shows that *CSHF* is a correct evaluator for *SHF*-processes w.r.t. the big-step semantics.

Instead of permitting to retry a transaction too often, the strategy from [HMPH08] can be applied: activate retried transactions only if some of the read TVars is modified. Our proofs can be used to show that this restriction of reductions is equivalent to the *CSHF*-semantics.

Note that omitting the send retry for abortion would make the implementation incorrect (non-adequate): a thread that runs into a loop if TVar x contains a 1 and returns otherwise, may block this thread indefinitely, which is not the case in the specification *SHF*.

5 Related Work

General remarks on STM are in [CBM⁺08,HLR10]. [CBM⁺08] argue that more research is needed to reduce the runtime overhead of STM. We believe that ease of maintenance of programs and the increased concurrency provided by STM may become more important in the future.

The paper which strongly influenced our work is [HMPH08]. *SHF* borrows from the operational semantics of STM Haskell, where differences are that our calculus is extended by futures, models also the call-by-need evaluation, but does not include exceptions and is restricted to monomorphic typing. However, since futures can be safely implemented in Concurrent Haskell by using `unsafeInterleaveIO` [SSS12], this difference is rather small. Moreover, our correctness proofs do not rely on having futures in the language and thus would also hold, if threads only have a thread identifier but not a resulting value.

describe the current implementation of STM Haskell in the Glasgow Haskell Compiler (GHC), however no formal treatment of this implementation is given. The approach taken in the GHC implementation is close to ours with the difference that instead of aborting transactions by the committing transaction (i.e. sending retries in *CSHF*), transactions abort and restart themselves by temporarily checking their local transaction log against the status of the global memory, and thus detecting conflicts. This is comparable to our approach, and we are convinced that the implementations are closely related. However, modelling their approach in a formal semantics would require more effort: For instance, in the GHC implementation comparing the local content (stored in the transaction log) and the global content of a TVar is done by pointer equality. Including this equality either has to be defined as a side condition of a reduction rule or has to be built in as a language construct. One also has to make design decisions like the following: Should indirections (i.e. bindings like $x = y$) be respected, by the pointer equality test? For example, if the local content of a TVar is x and the global content is y where the bindings are $x = z \mid z = s \mid y = w \mid w = z$, are x and y equal? Semantically, this is correct (but hard to express in a small-step semantics). Testing such cases in GHC's STM implementation shows that indirections are not respected in the current implementation.

In contrast the semantics of *CSHF* does not need to compare pointers, since it uses the registration mechanism. A difference between our semantics and STM-Haskell is that in *CSHF* reading the content of a TVar and thereafter writing the same content is always a modification which may result in aborting other transactions, while in STM-Haskell perhaps the pointer remains the same and thus no conflict modification is detected. Thus our approach may have slightly more restarts than the STM-Haskell implementation. A potential semantical problem in [HMPH08] is that exceptions may make local values of TVars visible outside the transaction.

A semantical investigation of STM is in [BBG09], where a call-by-name functional core language with concurrent processes is defined, and a contextual equivalence is used as equality. Also strong results are obtained by proving correctness of the monad laws and other program equivalences w.r.t. their semantics. However, [BBG09] only considers may-convergence in the contextual equivalence, which is too weak for reasoning about non-deterministic (and concurrent) calculi. Also, `seq` is missing in [BBG09] which is used in Haskell and known to change the semantics, such that validity of the monad laws only holds under further typing restrictions (see [SSS11]). A further difference to our work is that [BBG09] use pointer equality.

[HH09] propose to investigate correctness of an implementation, but stick to testing. [BT11] consider correctness of implementing STM in a small calculus with call-by-value reduction and a monadic extension similar to the STM/IO-extension in [HMPH08]. The main reasoning tool is looking for traces of effects, and arguing about commuting and shifting the effects within traces, where several important properties are proved. It is hard to compare the results with ours, however, from an abstract level, their proof method appears to ignore the should-convergence restriction: There is no argument on forced aborts of transactions.

6 Conclusion

We have presented a big-step semantics for STM-Haskell as a specification, and a small-step concurrent implementation. Using formal reasoning and the strong notion of contextual equivalence with may- and should-convergence we prove correctness of the implementation.

Further research directions are to consider smarter strategies for earlier aborts and retries of conflicting transactions, extending the language, for example by exceptions, and a polymorphic type system.

References

- BBG09. Johannes Borgström, Karthikeyan Bhargavan, and Andrew D. Gordon. A compositional theory for STM Haskell. In *Proc. Haskell '09*, pages 69–80. ACM, 2009.
- BT11. Annette Bieniusa and Peter Thiemann. Proving isolation properties for software transactional memory. In *Proc. ESOP'11*, volume 6602 of *LNCS*, pages 38–56, 2011.
- CBM⁺08. Calin Cascaval, Colin Blundell, Maged M. Michael, Harold W. Cain, Peng Wu, Stefanie Chiras, and Siddhartha Chatterjee. Software transactional memory: why is it only a research toy? *Commun. ACM*, 51(11):40–46, 2008.
- HH09. Liyang Hu and Graham Hutton. Towards a verified implementation of software transactional memory. In *Proc. TFP'08*, volume 9, pages 129–144. Intellect, 2009.
- HLR10. Tim Harris, James R. Larus, and Ravi Rajwar. *Transactional Memory, 2nd edition*. Synthesis Lectures on Computer Architecture. Morgan & Claypool Publishers, 2010.
- HMPH05. Tim Harris, Simon Marlow, Simon Peyton Jones, and Maurice Herlihy. Composable memory transactions. In *Proc. PPOPP'05*, pages 48–60. ACM, 2005.
- HMPH08. Tim Harris, Simon Marlow, Simon L. Peyton Jones, and Maurice Herlihy. Composable memory transactions. *Commun. ACM*, 51(8):91–100, 2008.
- RV07. Arend Rensink and Walter Vogler. Fair testing. *Inform. and Comput.*, 205(2):125–198, 2007.
- SSNSS08. Manfred Schmidt-Schauß, Joachim Niehren, Jan Schwinghammer, and David Sabel. Adequacy of compositional translations for observational semantics. In *Proc. IFIP TCS'08*, volume 273 of *IFIP*, pages 521–535. Springer, 2008.
- SSS08. David Sabel and Manfred Schmidt-Schauß. A call-by-need lambda-calculus with locally bottom-avoiding choice: Context lemma and correctness of transformations. *Math. Structures Comput. Sci.*, 18(03):501–553, 2008.
- SSS11. David Sabel and Manfred Schmidt-Schauß. A contextual semantics for Concurrent Haskell with futures. In *Proc. PPDP'11*, pages 101–112. ACM, 2011.
- SSS12. David Sabel and Manfred Schmidt-Schauß. Conservative concurrency in Haskell. In *Proc. LICS'12*, pages 561–570. IEEE, 2012.
- ST95. N. Shavit and D. Touitou. Software transactional memory. In *Proc. PODC'95*, pages 204–213. ACM, 1995.
- ST97. N. Shavit and D. Touitou. Software transactional memory. *Distributed Computing, Special Issue*, 10:99–116, 1997.
- vEPS90. Marko C. J. D. van Eekelen, Marinus J. Plasmeijer, and J. E. W. Smetsers. Parallel graph rewriting on loosely coupled machine architectures. In *Proc. CTRS'90*, volume 516 of *LNCS*, pages 354–369. Springer, 1990.